

Waller, W., Jain, A. (1978). On the monotonicity of the performance of Bayesian classifiers. *IEEE Transactions on Information Theory*, **IT-24**, 392-394.

Watanabe, S., *Pattern Recognition: Human and Mechanical*, Wiley and Sons, 1985.

Weiss, S.M., and Kulikowski, C. A. (1991). *Computer systems that learn*. Morgan Kauffman.

Wolpert, D. (1992). On the connection between in-sample testing and generalization error. *Complex Systems*, **6**, 47-94.

Wolpert, D. (1993). On overfitting avoidance as bias. SFI TR 93-03-016.

Wolpert, D. (1994a). The relationship between PAC, the Statistical Physics framework, the Bayesian framework, and the VC framework. In *The Mathematics of Generalization*. D. Wolpert (Ed.). Addison-Wesley.

Wolpert, D. (1994b). Filter likelihoods and exhaustive learning. In *Computational Learning Theory and Natural Learning Systems: Volume II*, S. Hanson et al. (Ed.'s). MIT Press.

Wolpert, D. (1995). On the Bayesian “Occam factors” argument for Occam’s razor. In *Computational Learning Theory and Natural Learning Systems: Volume III*, T. Petsche et al. (Ed.'s). MIT Press.

Wolpert, D. and Macready, W. (1995). No Free Lunch Theorems for Search. SFI TR 95-02-010. Submitted.

Wolpert, M, Grossman, T., and Knill, E. (1995). Off-Training-Set Error for the Gibbs and the Bayes Optimal Generalizers. Submitted.

- Bridle, J. (1989). Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition. In F. Fougelman-Soulie and J. Herault (Eds.), *Neuro-computing: Algorithms, architectures, and applications*. Springer-Verlag.
- Dietterich, T. (1990). Machine Learning. *Annu. Rev. Comput. Sci.*, **4**, 255-306.
- Drucker, H. et al. (1993). Improving performance in neural networks using a boosting algorithm. In *Neural Information Processing Systems 5*, S. Hanson et al. (Eds). Morgan-Kauffman.
- Duda, R., and Hart, P. (1973). *Pattern classification and scene analysis*. Wiley and Sons.
- Hughes, G. (1968). On the mean accuracy of statistical pattern recognizers. *IEEE Transactions on Information Theory*, **IT-14**, 55-63.
- Kearns, M. J., et al. 'Towards efficient agnostic learning, in *Proceedings of the 5th annual workshop on Computational Learning Theory*, ACM Press, NY, NY, 1992.
- Mitchell, T. (1982). Generalization as search, *Artificial Intelligence*, **18**, 203-226.
- Mitchell T., Blum, A. (1994). *Course notes for Machine Learning*, CMU.
- Murphy, P., Pazzani, M. (1994). Exploring the decision forest: an empirical investigation of Occam's razor in decision tree induction. *Journal of Artificial Intelligence Research*, **1**, 257-275.
- Natarajan, B. (1991). *Machine Learning: A theoretical approach*. Morgan Kauffman, San Mateo, CA.
- Perrone, M. (1993). Improving regression estimation: averaging methods for variance reduction with extensions to general convex measure optimization. Ph.D. thesis, Brown Univ. Physics Dept.
- Plutowski, M., et al. (1994). Cross-validation estimates integrated mean squared error. In *Advances in neural information processing systems 6*, Cowan et al. (Ed.'s), Morgan Kauffman, CA.
- Schaffer, C., (1993). Overfitting avoidance as bias. *Machine Learning*, **10**, 153-178.
- Schaffer, C. (1994). A conservation law for generalization performance. In Cohen and Hirsh (Ed.'s), *Machine Learning: Proceedings of the Eleventh International Conference*. Morgan Kauffman, San Francisco.
- Schapire, R. (1990). The strength of weak learnability. *Machine Learning*. **5**, 197-227.
- Vapnik, V. (1982). Estimation of dependences based on empirical data. Springer-Verlag.
- Vapnik, V., and Bottou, L. (1993). Local algorithms for pattern recognition and dependencies estimation. *Neural Computation*, **5**, 893-909.

is ordered and perhaps has repeats, the value of $P(d_{\mathbf{X}} \mid \text{no punt}, \phi, m)$ when all the $d_{\mathbf{X}}(i) \leq k$ is k^{-m} . Similarly, $P(d_{\mathbf{X}} \mid \text{punt}, \phi, m) = 0$ unless at least one of the $d_{\mathbf{X}}(i) > k$, and when it's non-zero it equals some constant set by k and m .

It's also true that $E(C_{\text{OTS}} \mid \phi, d_{\mathbf{X}})$ is not drastically different if one considers $d_{\mathbf{X}}$'s with a different m' . Accordingly, our summand doesn't vary drastically between $d_{\mathbf{X}}$'s of one m' and $d_{\mathbf{X}}$'s of another. Since $n \gg m$ and $\pi(x)$ is uniform though, almost all of the terms in the sum have $m' = m$. Pulling this all together, we see that to an arbitrarily good approximation (for large enough n relative to m), we can take $m' = m$. So (E.1) becomes

$$\text{E.2) } E(C_{\text{OTS}} \mid \phi, (\text{no punt}), m) = \sum_{d_{\mathbf{X}}} E(C_{\text{OTS}} \mid \phi, d_{\mathbf{X}}) P(d_{\mathbf{X}} \mid (\text{no punt}), m' = m).$$

Now consider conditioning on 'no punt', in which case all the $d_{\mathbf{X}}(i) \leq k$. For such a situation, for $m' = m$, $E(C_{\text{OTS}} \mid \phi, d_{\mathbf{X}}) = (n - k) / (n - m)$. In contrast, consider having a punt signal, in which case at least one $d_{\mathbf{X}}(i) > k$. Now $E(C_{\text{OTS}} \mid \phi, d_{\mathbf{X}}) \leq (n - k - 1) / (n - m) < (n - k) / (n - m)$.

Combining this with (E.2), we get $E(C_{\text{OTS}} \mid \phi, \text{punt}, m) < E(C_{\text{OTS}} \mid \phi, \text{no punt}, m)$. QED.

References

- Anthony M. and Biggs N. (1992). *Computational Learning Theory*. Cambridge University Press.
- Berger, J. (1985). *Statistical decision theory and Bayesian analysis*. Springer-Verlag.
- Bernardo, J. Smith, A. (1994). *Bayesian Theory*. Wiley and Sons, NY.
- Berger, J., and Jeffreys, W. (1992). Ockham's razor and Bayesian analysis. *American Scientist*, **80**, 64-72.
- Blumer, A., et alia (1987). Occam's razor. *Information Processing Letters*, **24**, 377-380.
- Blumer, A., et alia (1989). Learnability and Vapnik-Chervonenkis dimension. *Journal of the ACM*, **36**, 929-965.

ϕ_2 are a pair of discrete-valued vectors, and α is a real-valued vector indexed by a value for ϕ_1 and one for ϕ_2 ; transform α so that its dependence on ϕ_2 is rearranged in some arbitrary - though invertible - fashion. Performing this transformation is equivalent to mapping the space of all ϕ_2 vectors into itself in a one-to-one manner. The Jacobian of this transformation is 1, and the transformation doesn't change the functional form of the constraint forcing α to lie on a simplex. (I.e., $\sum_{\phi_1\phi_2} \alpha_{\phi_1\phi_2} = 1$ and for all $\phi_1\phi_2$, $\alpha_{\phi_1\phi_2} \geq 0$, where double-prime indicates the new ϕ_2 indices.) So expressed in this new coordinate system, the integral is $\int d\alpha \{ \alpha_{\phi_1, \phi_2^*} / \sum_{\phi_1} \alpha_{\phi_1} P(d | \phi_1) \}$, where ϕ_2^* is a new index corresponding to the old index ϕ_2 . Since this integral must have the same value as our original integral, and since ϕ_2^* is arbitrary, we see that that integral is independent of ϕ_2 , and therefore can only depend on the values of d and ϕ_1 .

This means that we can rewrite our sum over all ϕ as

$$\sum_{\phi_1\phi_2} P(c | \phi_2, d) P(d | \phi_1) \text{func}_1\{\phi_1, d\}$$

for some function “ $\text{func}_1(\cdot)$ ”. In other words, the α -average of $P(c | d, \alpha)$ is proportional to $\sum_{\phi_2} P(c | \phi_2, d)$, where the proportionality constant depends on d . Since $P(c | \phi, d) = P(c | \phi_2, d)$ (see above), our sum is proportional to $\sum_{\phi_1\phi_2} P(c | \phi, d) = \sum_{\phi} P(c | \phi, d)$. By theorem (4), this sum equals $\Lambda(c) / r$.

So the uniform α -average of $P(c | d, \alpha) = \text{func}_2(d) \Lambda(c) / r$ for some function “ $\text{func}_2(\cdot)$ ”. Since $\sum_c P(c | d, \alpha) = 1$, the sum over C values of the uniform α -average of $P(c | d, \alpha)$ must be independent of d (it must equal 1). Therefore $\text{func}_2(d)$ is independent of d . Since we know that $\Lambda(c) / r$ is properly normalized over c , we see that $\text{func}_2(d)$ in fact equals 1. QED.

APPENDIX E. PROOF OF THEOREM (10)

First use the fact that given ϕ , $d_{\mathbf{X}}$ determines whether there is a punt signal, to write

$$E.1) E(C_{\text{OTS}} | \phi, (\text{no}) \text{ punt}, m) = \sum_{d_{\mathbf{X}}} E(C_{\text{OTS}} | \phi, d_{\mathbf{X}}) P(d_{\mathbf{X}} | (\text{no}) \text{ punt}, \phi, m)$$

Next, without loss of generality, let the x 's for which $\phi(x) = 0$ be $1, \dots, k$, so that $\phi(x) = 1$ for $x = k + 1, \dots, n$. Then $P(d_{\mathbf{X}} | \text{no punt}, \phi, m) = 0$ unless all the $d_{\mathbf{X}}(i) \leq k$. Since $\pi(x)$ is uniform, and d

empirical error s .

Proof that active learning has a vertical likelihood: Let d_k refer to the first k input-output pairs in the training set d , and $d(i)$ to the i 'th such pair. Then $P(d_m | f) = P(d(m) | f, d_{m-1}) P(d_{m-1} | f) = P(d_Y(m) | d_X(m), f, d_{m-1}) P(d_X(m) | f, d_{m-1}) P(d_{m-1} | f)$. By hypothesis, in active learning $P(d_X(m) | f, d_{m-1}) = P(d_X(m) | d_{m-1})$. So long as it is also true that $P(d_Y(m) | d_X(m), f, d_{m-1}) = P(d_Y(m) | d_X(m), f)$ is independent of $f(x \neq d_X(m))$, by induction we have a vertical likelihood.

APPENDIX D. PROOF OF THEOREM 8

The task before us is to calculate the average over all α of $P(c | d, \alpha)$. To that end, write the average as (proportional to) $\int d\alpha [\sum_{\phi} P(\phi | d, \alpha) P(c | \phi, d, \alpha)]$, where as usual the integral is restricted to the r^n -dimensional simplex. Rewrite this integral as $\int d\alpha [\sum_{\phi} P(\phi | \alpha) P(c | \phi, d, \alpha) P(d | \phi, \alpha) / P(d | \alpha)] = \int d\alpha [\sum_{\phi} \alpha_{\phi} P(c | \phi, d) P(d | \phi)] / [\sum_{\phi'} \alpha_{\phi'} P(d | \phi')]$, where ϕ' is a dummy ϕ value. Rewrite this in turn as

$$\sum_{\phi} P(c | \phi, d) P(d | \phi) \left\{ \int d\alpha \frac{\alpha_{\phi}}{\sum_{\phi'} \alpha_{\phi'} P(d | \phi')} \right\}.$$

As in the proof of theorem (2), break up ϕ into two components, ϕ_1 and ϕ_2 , where ϕ_1 fixes the values of ϕ over the \mathbf{X} values lying inside d_X , and ϕ_2 fixes it over the values outside of d_X . We must find how the terms in our sum depend on ϕ_1 and ϕ_2 .

First, write $P(c | \phi, d) = \sum_H P(h | d) P(c | h, \phi, d)$. By definition, for OTS error $P(c | h, \phi, d)$ is independent of ϕ_1 . This allows us to write $P(c | \phi, d) = P(c | \phi_2, d)$.

Next, since we're restricting attention to vertical likelihoods, $P(d | \phi)$ only depends on ϕ_1 . So we can write the term in the curly brackets as $\int d\alpha [\alpha_{\phi_1 \phi_2} / \sum_{\phi_1' \phi_2'} \alpha_{\phi_1' \phi_2'} P(d | \phi_1')] = \int d\alpha [\alpha_{\phi_1 \phi_2} / \sum_{\phi_1'} \alpha_{\phi_1'} P(d | \phi_1')]$ with obvious notation. Since we're assuming that for no ϕ does $P(d | \phi)$ equal zero exactly, the denominator sum is always non-zero.

Now change variables in the integral over α by rearranging the ϕ_2 indices of α . In other words: ϕ_1 and

tionality constant. By normalization, the resultant value of our constant must be the reciprocal of $\int df(x \notin d_{\mathbf{X}}) 1$. QED.

Proof of theorem (7): $\int d\alpha P(c | m, \alpha) = \int d\alpha [\sum_{\phi} P(\phi | m, \alpha) P(c | m, \alpha, \phi)]$, where the integral is restricted to the r^n -dimensional simplex. This can be rewritten as $\int d\alpha [\sum_{\phi} \alpha_{\phi} P(c | \phi, m, \alpha)]$, since we assume that the probability of ϕ has nothing to do with the number of elements in d . Similarly, once ϕ is fixed, the probability that $C = c$ doesn't depend on α , so our average equals $\int d\alpha [\sum_{\phi} \alpha_{\phi} P(c | \phi, m)]$. Write this as $\sum_{\phi} P(c | \phi, m) [\int d\alpha \alpha_{\phi}]$. By symmetry, the term inside the square brackets is independent of ϕ . Therefore the average over all $P(\phi)$ of $P(c | m)$ is proportional to $\sum_{\phi} P(c | \phi, m)$. Using theorem (5) and normalization, this establishes theorem (7). QED.

Proof of corollary (3): Follow along with the proof of theorem (7). Instead of $\int d\alpha \alpha_{\phi}$, we have $\int d\alpha G(\alpha) \alpha_{\phi}$. (For present purposes, the delta and Heaviside functions that force α to stay on the unit simplex are implicit.) By assumption, $G(\alpha)$ is unchanged under the bijection of replacing all vectors α_i with new vectors identical to the old, except that the components for $i = \phi$ and $i = \phi'$ are interchanged. This is true for all ϕ' and ϕ . Accordingly, our integral is independent of ϕ , which suffices to prove the result. QED.

Proof of theorem (9): To evaluate $P(c | s, d)$ for uniform $P(f)$, write it as $\int df P(c | s, d, f) P(f | d, s)$. Next write $P(f | d, s) = P(s | f, d) P(f | d) / P(s | d)$. Note though that $P(s | f, d) = P(s | d)$ (see beginning of section (4)), and recall that we are implicitly assuming that $P(s | d) \neq 0$. So we get $P(c | s, d) = \int df P(c | s, d, f) P(d | f)$, up to an overall d -dependent proportionality constant. Now proceed as in the proof of theorem (2) by breaking the integral into two integrals, one over $f(x \in d_{\mathbf{X}})$, and one over $f(x \notin d_{\mathbf{X}})$. The result is $P(c | s, d) = \Lambda(c) / r$, up to an overall d -dependent proportionality constant. By normalization, that constant must equal 1. This establishes theorem (9). QED.

Example of non-NFL behavior of s -conditioned distributions: Let $P(h | d) = \delta(h, h^*)$ for some h^* , let $\pi(x)$ be uniform, use zero-one loss, assume a noise-free IID likelihood, and take $m = 1$. Then we can write $E(C_{OTS} | s, f, m = 1) = E(C_{OTS} | s, f, m = 1, h = h^*) = [n C'_{IID}(f, h^*) - s] / (n - 1)$. (Note that C'_{IID} is independent of d , and that for zero-one loss $n \times C'_{IID}(f, h^*)$ is the number of disagreements between h^* and f over all of \mathbf{X} .) No matter what f is, this grows as s shrinks. Since C_{OTS} can only have two values, this means that as s grows, $P(C_{OTS} | f, s, m = 1)$ gets biased towards the lower of the two possible values of c_{OTS} . So we do not have NFL behavior for the uniform average over f of $P(c | f, s, m)$ - that average depends on the

Proof of the claim concerning the “random learning algorithm”, made just below lemma (1): By lemma 1, $P(c | f, d) = \sum_{y_H, y_F, q} \delta(c, L(y_H, y_F)) P(y_H | q, d) P(y_F | q, f) P(q | d)$. However for OTS error $q \notin d_X$, and therefore for the random learning algorithm, for all q and d in our sum, $P(y_H | q, d) = 1 / r$, independent of y_H (recall that there are r elements in Y). If we have a symmetric homogenous loss function, this means that we can replace $\sum_{y_H} \delta(c, L(y_H, y_F)) P(y_H | q, d)$ with $\Lambda(c) / r$. Since this is independent of d and f , $P(c | d) = \Lambda(c) / r$ for all training sets d , as claimed. QED.

Proof of the “implication of lemma (1)”, made just below lemma (1): Uniformly average the expression for $P(c | f, d)$ in lemma (1) over all targets f . The only place f occurs in the sum in lemma (1) is in the third term, $P(y_F | q, f)$. Therefore our average replaces that third term with some function $\text{func}(y_F, q)$. By symmetry though, the uniform f -average of that third term in the sum must be the same for all test set inputs q and outputs y_F . Accordingly $\text{func}(y_F, q)$ is some constant. Now the sum over y_F of this constant must equal 1 (to evaluate that sum of the f -average of $P(y_F | q, f)$, interchange the sum over y_F with the average over f). Therefore our constant must equal $1 / r$. The implication claimed is now immediate. QED.

Proof of theorem (2): We can replace the sum over all q that gives $P(c | f, d)$ (lemma (1)) with a sum over those q lying outside of d_X . Accordingly, for such a $P(q | d)$, $P(c | f, d)$ is independent of the values of $f(x \in d_X)$. (For clarity, the second argument of f is being temporarily suppressed.)

Noting that $P(d | f)$ is vertical, next average both sides of our equation for $P(c | f, m)$ uniformly over all f and pull the f -average inside the sum over d . Since $P(c | f, d)$ and $P(d | f)$ depend on separate parts of f (namely $f(x \notin d_X)$ and $f(x \in d_X)$ respectively), we can break the average over f into two successive averages, one operating on each part of f , and thereby get

$$\sum_d \left[\int df(x \notin d_X) P(c | f, d) \int df(x \in d_X) P(d | f) \right] / \int df(x \notin d_X) \int df(x \in d_X) 1 .$$

But since $P(c | f, d)$ is independent of the values of $f(x \in d_X)$, uniformly averaging it over all $f(x \in d_X)$ is equivalent to uniformly averaging it over all f . By theorem (1), such an average is independent of d . Therefore we can pull that average out of the sum over d , and get theorem (2). QED.

Proof of theorem (3): Write $P(c | d)$ for a uniform $P(f) \propto \int df P(c | d, f) P(d | f)$, where the proportionality constant depends on d . Break up the integral over f into an integral over $f(x \in d_X)$ and one over $f(x \notin d_X)$, exactly as in the proof of theorem (2). Absorb $\int df(x \in d_X) P(d | f)$ into the overall (d -dependent) propor-

$$\delta \equiv \max_{\{d_{\mathbf{X}}, f, h\}} \{ C'_{\text{OTS}}(d, f, h) [1 - \pi(\mathbf{X} - d_{\mathbf{X}})] + C'_{\text{TS}}(d, f, h) \pi(d_{\mathbf{X}}) \}.$$

Then write

$$\begin{aligned} P_{C'_{\text{IID}} | Z} (c'_{\text{IID}} + \kappa | z) \\ = \sum_d \int df dh \delta(c'_{\text{IID}} - C'_{\text{OTS}}(d, f, h)) \times P(d, f, h | z) \end{aligned}$$

for some constant κ where $|\kappa| \leq \delta$.

Now for non-pathological z , $P(c'_{\text{IID}} | z)$ is a continuous function of c'_{IID} . (Recall in particular that in this paper, no event has exactly zero probability; see appendix A.) So for such a z , for ε sufficiently small, $\delta \rightarrow 0$ and therefore $\kappa \rightarrow 0$, and we can approximate

$$\begin{aligned} P(c'_{\text{IID}} | z) \\ = \sum_d \int df dh \delta(c'_{\text{IID}} - C'_{\text{OTS}}(d, f, h)) \times P(d, f, h | z). \end{aligned}$$

But this just equals $P_{C'_{\text{OTS}} | F, M}(c'_{\text{IID}} | z)$. So for $n \gg r$ and non-pathological z and $\pi(\cdot)$, the distribution over c'_{IID} is the same as that over c'_{OTS} . QED.

APPENDIX C. MISCELLANEOUS PROOFS

For clarity of the exposition, several of the more straight-forward proofs in the paper are collected in this appendix.

Proof of lemma 1. Write $P(c | f, d) = \sum_{y_H, y_F, q} P(c | y_H, y_F, q, f, d) P(y_H | y_F, q, f, d) P(y_F, q | f, d)$. Re-writing the summand, we get $P(c | f, d) = \sum_{y_H, y_F, q} \delta(c, L(y_H, y_F)) P(y_H | y_F, f, q, d) P(y_F, q | f, d)$.

Now $P(y_H | y_F, f, q, d) = \int dh P(y_H | y_F, h, f, q, d) P(h | y_F, f, q, d) = \int dh P(y_H | h, q, d) P(h | q, d)$ (see point (11) in the EBF section). This just equals $P(y_H | q, d)$. (However it is *not* true in general that $P(y_H | y_F, d) = P(y_H | d)$. See [Wolpert et al. 1995].) Plugging in gives the result. QED.

ately carry over to IID error. This is proven formally in Appendix B.

APPENDIX B. PROOF THAT DISTRIBUTIONS OVER C'_{IID} EQUAL THOSE OVER C'_{OTS} WHENEVER $n \gg r$, FOR NON-PATHOLOGICAL $\pi(\cdot)$

To prove the assertion, with slight abuse of terminology write $C'_{\text{IID}} = C'_{\text{OTS}} \pi(\mathbf{X} - \mathbf{d}_{\mathbf{X}}) + C'_{\text{TS}} \pi(\mathbf{d}_{\mathbf{X}})$, where “TS” means error on the training set, defined in the obvious way, and $\pi(A) \equiv \sum_{\mathbf{x} \in A} \pi(\mathbf{x})$ (see [Wolpert et al., 1995]). Then for any set of one or more random variables Z taking value z , we have

$$\begin{aligned}
 P(c'_{\text{IID}} | z) &= \sum_{\mathbf{d}} \int dc'_{\text{OTS}} dc'_{\text{TS}} P(c'_{\text{IID}} | c'_{\text{OTS}}, c'_{\text{TS}}, \mathbf{d}, z) \times P(c'_{\text{OTS}}, c'_{\text{TS}}, \mathbf{d} | z) \\
 &= \sum_{\mathbf{d}} \int dc'_{\text{OTS}} dc'_{\text{TS}} \delta(c'_{\text{IID}} - c'_{\text{OTS}} \pi(\mathbf{X} - \mathbf{d}_{\mathbf{X}}) - c'_{\text{TS}} \pi(\mathbf{d}_{\mathbf{X}})) \times P(c'_{\text{OTS}}, c'_{\text{TS}}, \mathbf{d} | z) \\
 &= \sum_{\mathbf{d}} \int dc'_{\text{OTS}} dc'_{\text{TS}} \delta(c'_{\text{IID}} - c'_{\text{OTS}} \pi(\mathbf{X} - \mathbf{d}_{\mathbf{X}}) - c'_{\text{TS}} \pi(\mathbf{d}_{\mathbf{X}})) \times \\
 &\quad \int df dh P(c'_{\text{OTS}}, c'_{\text{TS}} | \mathbf{d}, f, h, z, m) P(\mathbf{d}, f, h | z).
 \end{aligned}$$

Now again abuse terminology slightly and write

$$P(c'_{\text{OTS}}, c'_{\text{TS}} | \mathbf{d}, f, h, z) = \delta(c'_{\text{OTS}} - C'_{\text{OTS}}(\mathbf{d}, f, h)) \times \delta(c'_{\text{TS}} - C'_{\text{TS}}(\mathbf{d}, f, h)),$$

where the statistical dependencies of C'_{OTS} and C'_{TS} are made explicit by writing them as functions. Plugging in we get

$$\begin{aligned}
 P(c'_{\text{IID}} | z) &= \sum_{\mathbf{d}} \int df dh \delta(c'_{\text{IID}} - C'_{\text{OTS}}(\mathbf{d}, f, h) \pi(\mathbf{X} - \mathbf{d}_{\mathbf{X}}) - C'_{\text{TS}}(\mathbf{d}, f, h) \pi(\mathbf{d}_{\mathbf{X}})) \times P(\mathbf{d}, f, h | z).
 \end{aligned}$$

Define $\varepsilon \equiv \max_{\mathbf{d}_{\mathbf{X}}} \pi(\mathbf{d}_{\mathbf{X}})$, so $\min_{\mathbf{d}_{\mathbf{X}}} \pi(\mathbf{X} - \mathbf{d}_{\mathbf{X}}) = 1 - \varepsilon$. Now whenever $n \gg r$, so long as there are no sharp peaks in $\pi(\cdot)$, $\varepsilon \rightarrow 0$. However because a delta function is not a continuous function, taking the limit as $\varepsilon \rightarrow 0$ of our expression for $P(c'_{\text{IID}} | z)$ is not immediately equivalent to setting the $\pi(\mathbf{X} - \mathbf{d}_{\mathbf{X}})$ and $\pi(\mathbf{d}_{\mathbf{X}})$ inside the delta function to 1 and to 0 respectively. We can circumvent this difficulty rather easily though. To do that, define

$$E(C \mid f, h, d) = \sum_{y_H, y_F, q} L(y_H, y_F) h(q, y_H) f(q, y_F) P(q \mid d).$$

In much of supervised learning, an expression like that on the right-hand side of this equation is called the “generalization error”. In other words, instead of the error C used here, in much of supervised learning one uses an alternative error C' , defined by $C'(f, h, d) \equiv E(C \mid f, h, d)$, i.e., $P(c' \mid f, h, d) = \delta(c', E(C \mid f, h, d))$.

- Note that in general, the set of allowed values of C is not the same as the set of allowed values of C' . In addition, distributions over C do not set those over C' . For example, knowing $P(c \mid d)$ need not give $P(c' \mid d)$ or vice-versa.⁶ However many branches of supervised learning theory (e.g., much of computational learning theory) are concerned with quantities of the form “ $P(\text{error} > \epsilon \mid \dots)$ ”.⁷ For such quantities, whether one takes “error” to mean C or (as is conventional) C' may change the results, and in general one cannot directly deduce the result for C from that for C' (or vice versa).

Where appropriate, subscripts $_{OTS}$ or $_{IID}$ on c' will indicate which kind of $P(q \mid d)$ is being used.

- Fortunately, most of the results derived in this paper apply equally well to both probabilities of C and probabilities of C' . For reasons of space though, I will only work out the results explicitly for C . However note that we can immediately equate expectations of C that are not conditioned on q , y_H or y_F with the same expectations of C' . For example,

$$E(C \mid d) = \int dhdf E(C \mid f, h, d) P(f, h \mid d) = \int dhdf C'(f, h, d) P(f, h \mid d) = \int dhdf E(C' \mid f, h, d) P(f, h \mid d) = E(C' \mid d).$$

So when cast in terms of expectation values, any (appropriately conditioned) results automatically apply to C' as well as C .

MISCELLANEOUS

- For most purposes, it is implicitly assumed that no probabilities equal zero *exactly* (although some probabilities might be infinitesimal). That way we have never have to worry about dividing by probabilities, and in particular never have to worry about whether conditional probabilities are well-defined. So as an example, phrases like “noise-free” are taken to mean infinitesimal noise rather than exactly zero noise.
- It is important to note that in general, for non-pathological $\pi(\cdot)$, in the limit where $n \gg r$, distributions over c'_{IID} are identical to distributions over c'_{OTS} . In this sense theorems concerning OTS error immedi-

rithm only sees d in making its guess, not f). This means that $P(h, f | d) = P(h | d) P(f | d)$, and therefore $P(f | h, d) = P(f | d)$.

As an example of the importance of this assumption, note that it implies that $P(y_F | y_H, d, q) = P(y_F | d, q)$:

Proof: Expand $P(y_F | y_H, d, q) = \int df f(q, y_F) P(f | d, q) \int dh h(q, y_H) P(h | f, d, q)$. Since $P(h | f, d, q) = P(h | d)$, this integral is proportional to $\int df f(q, y_F) P(f | d, q)$, where the proportionality constant depends on d and q . However $\int df f(q, y_F) P(f | d, q) = P(y_F | d, q)$. Due to normalization, this means that the proportionality constant equals 1, and we have established the proposition. QED.

Our assumption does not imply that $P(y_F | y_H, d) = P(y_F | d)$ however. Intuitively, for a fixed learning algorithm, knowing y_H and d tells you something about q , and therefore (in conjunction with knowledge of d) something about y_F , that d alone does not.

- The “posterior” is the Bayesian inverse of the likelihood, $P(f | d)$. The phrase “the prior” usually refers to $P(f)$.
- Some schemes can be cast into this framework in more than one way. As an example, consider softmax [Bridle 1989], where each output neuron indicates a different possible event, and the real values the neurons take in response to an input are interpreted as input-conditioned probabilities of the associated events. For this scheme one could either i) take \mathbf{Y} to be the set of “possible events”, so that the h produced by the algorithm is not single-valued, or ii) take \mathbf{Y} to be (the computer’s discretization of) the real-valued vectors that the set of output neurons can take on, in which case h is single-valued, and \mathbf{Y} itself is interpreted as a space of probability distributions. Ultimately, which interpretation one adopts is determined by the relationship between C and H . (Such relationships are discussed below.)

“GENERALIZATION ERROR”

- Note that $E(C | f, h, d) = \sum_{y_H, y_F, q} E(C | f, h, d, y_H, y_F, q) P(y_H, y_F, q | f, h, d)$. Due to our definition of C , the first term in the sum equals $L(y_H, y_F)$. The second term equals $P(y_H | h, q, f, d, y_F) P(y_F | q, f, d, h) P(q | d, f, h)$. This in turn equals $h(q, y_H) f(q, y_F) P(q | f, h, d)$. In addition, $P(q | f, h, d) = P(q | d)$ always in this paper. Therefore

cordingly, all the usual rules concerning f apply as well to ϕ . (For example, $P(h | d, \phi) = P(h | d)$.) When I wish to make clear what ϕ sets f , I will write f_ϕ , as above; ϕ simply serves as an index on f . (In general, depending on $N(\cdot)$, it might be that more than one ϕ labels the same f , but this won't be important for the current analysis.) So when I say something like “vertical $P(d | \phi)$ ” it is implicitly understood that I mean vertical $P(d | f_\phi)$.

- When I say that I am “only allowing” these kinds of f , I will mean that whenever ‘ f ’ is written, it is assumed to be related to a ϕ in this manner - all other f implicitly have an infinitesimal prior probability.
- Note that the $N(\cdot)$ introduced here is the noise process operating in the generation of the test set, and need not be the same as the noise process in the generation of the training set. As an example, it is common in the neural net literature to generate the training set by adding noise to a single-valued function from \mathbf{X} to \mathbf{Y} , $\phi(\cdot)$, but to measure error by how well the resulting h matches that underlying $\phi(\cdot)$, not by how well Y_H values sampled from h match Y_F values formed by adding noise to $\phi(\cdot)$. In the ϕ - N terminology, this would mean that although $P(d | f)$ may be formed by corrupting some function $\phi(\cdot)$ with noise (in either \mathbf{X} and/or \mathbf{Y}), $P(y_F | f, q)$, which is used to measure test set error, is determined by a noise-free $N(\cdot)$, $N(y_F | q, \phi(q)) = \delta(y_F, \phi(q))$.
- Of special importance will be those noise-processes for which for each q , the uniform ϕ -average of $P(y_F | q, \phi)$ is independent of y_F . (Note this doesn't exclude q -dependent noise processes). I will call such a (test set) noise process “homogenous”. Intuitively, such noise processes have no *a priori* preference for one \mathbf{Y} value over another. As examples, the noise-free testing mentioned just above is homogenous, as is a noise process that when it takes in a value of $\phi(q)$, produces the same value with probability z and all other values with (identical) probabilities $(1 - z) / (r - 1)$.

COUPLING ALL THIS TO SUPERVISED LEARNING

- Any (!) learning algorithm (or “generalizer”) is simply a distribution $P(h | d)$. It is “deterministic” if the same d always gives the same h (i.e., if for fixed d $P(h | d)$ is a delta function about one particular h).
- There are many equalities that are assumed in supervised learning, but that do not merit explicit delineation. For example, it is implicitly assumed that $P(h | q, d) = P(h | d)$, and therefore that $P(q | d, h) = P(q | d)$.
- One assumption that does merit explicit delineation is that $P(h | f, d) = P(h | d)$ (i.e., the learning algo-

iterating the following “independent identically distributed” (IID) procedure: Choose \mathbf{X} values according to a “sampling distribution” $\pi(\mathbf{x})$, and then sample f at the those points to get associated \mathbf{Y} values.⁴ More formally, this very common scheme is equivalent to the following “likelihood”, presented previously as equation (2.1):

$$\begin{aligned} P(d | f) &= P(d_{\mathbf{Y}} | f, d_{\mathbf{X}}) P(d_{\mathbf{X}} | f) = P(d_{\mathbf{Y}} | f, d_{\mathbf{X}}) P(d_{\mathbf{X}}) \text{ (by assumption)} \\ &= \prod_{i=1}^m [\pi(d_{\mathbf{X}}(i)) f(d_{\mathbf{X}}(i), d_{\mathbf{Y}}(i))]. \end{aligned}$$

There is no *a priori* reason for $P(d | f)$ to have this form however. For example, in “active learning” or “query-based” learning, successive values (as i increases) of $d_{\mathbf{X}}(i)$ are determined by the preceding values of $d_{\mathbf{X}}(i)$ and $d_{\mathbf{Y}}(i)$. As another example, typically $P(d | f)$ will not obey equation (2.1) if testing and training are not governed by the same $P(y | x)$. (Recall that f governs the generation of test sets.) To see this, let t be the random variable $P(y | x)$ governing the generation of training sets. Then $P(d | f) = \int dt P(d | t) P(t | f)$. Even if $P(d | t) = \prod_{i=1}^m [\pi(d_{\mathbf{X}}(i)) t(d_{\mathbf{X}}(i), d_{\mathbf{Y}}(i))]$, unless $P(t | f)$ is a delta function about $t = f$, $P(d | f)$ need have the form specified in equation (2.1).

- I will say that $P(d | f)$ is “vertical” if it is independent of the values of $f(x \notin d_{\mathbf{X}})$. Any likelihood of the form given in equation (2.1) is vertical, by inspection. In addition, as discussed in section 4, active learning usually has a vertical likelihood. However some scenarios in which $t \neq f$ do not have vertical likelihoods.⁵
- In the case of “IID error” (the conventional error measure), $P(q | d) = \pi(q)$. In the case of OTS error, $P(q | d) = [\delta(q \notin d_{\mathbf{X}}) \pi(q)] / [\sum_q \delta(q \notin d_{\mathbf{X}}) \pi(q)]$, where $\delta(z) \equiv 1$ if z is true, 0 otherwise. Strictly speaking, OTS error is not defined when $m' = n$.

Where appropriate, subscripts _{OTS} or _{IID} on c will indicate which kind of $P(q | d)$ is being used.

FUNCTION+NOISE TARGETS

- In this paper I will consider in some detail those cases where we only allow those f that can be viewed as some single-valued function ϕ taking \mathbf{X} to \mathbf{Y} with a fixed noise process in \mathbf{Y} superimposed.⁸ To do this, I will (perhaps only implicitly) fix a noise function N that is a probability distribution over \mathbf{Y} , conditioned on $\mathbf{X} \times \mathbf{Y}$; N is a probability distribution over y_F , conditioned on the values of q and $\phi(q)$. (Note that there are r^n such functions $\phi(\cdot)$.)

Given $N(\cdot)$, each ϕ specifies a unique f via $P(y_F | f_\phi, q) = f_\phi(q, y_F) = N(y_F | q, \phi(q)) = P(y_F | q, \phi)$. Ac-

y_F).

All of this is formalized as follows.

- The F random variable parameterizes the Q -conditioned distribution over Y_F : $P(y_F | f, q) = f(q, y_F)$. In other words, f determines how test set elements y_F are generated for a test set point q . So Y_F and Q are the random variables whose relationship to F allows F to be intuitively viewed as an “ \mathbf{X} -conditioned distribution over \mathbf{Y} ” - see above.
- The variable Y_H meets similar requirements: $P(y_H | h, q) = h(q, y_H)$, and this relationship between Y_H , Q , and H is what allows one to view H as intuitively equivalent to an “ \mathbf{X} -conditioned distribution over \mathbf{Y} ”.
- For the purposes of this paper, the random variable cost C is defined by $C = L(Y_H, Y_F)$, where $L(., .)$ is called a “loss function”. As examples, zero-one loss has $L(a, b) = 1 - \delta(a, b)$, where $\delta(a, b)$ is the Kronecker delta function, and quadratic loss has $L(a, b) = (a - b)^2$. (Zero-one loss is assumed in almost all of computational learning theory.)

It is important to note though that in general C need not correspond to such a loss function. For example, “logarithmic scoring” has $c = -\sum_y f(q, y) \ln[h(q, y)]$, and does not correspond to any $L(y_F, y_H)$.

- For many L ’s the sum over y_F of $\delta(c, L(y_H, y_F))$ is some function $\Lambda(c)$, independent of y_H . I will call such L ’s “homogenous”. Intuitively, such L ’s have no *a priori* preference for one \mathbf{Y} value over another. As examples, the zero-one loss is homogenous. So is the squared difference between y_F and y_H if they are viewed as angles, $L(y_F, y_H) = [(y_F - y_H) \bmod \pi]^2$.

Note that one can talk of an L ’s being homogenous for certain values of c . For example, the quadratic loss isn’t homogenous over all c , but it is for $c = 0$. The results presented in this paper that rely on homogeneity of L usually hold for a particular c so long as L is homogenous for that c , even if L isn’t homogenous for all c .

THE RELATIONSHIP BETWEEN F , D AND Q

- Note that f is a distribution governing test set data (it governs the outputs associated with q), and in general it need not be the same as the distribution governing training set data. Unless explicitly stated otherwise though, I will assume that both training sets and test sets are generated via f .
- Often when training and testing sets are generated by the same $P(y | x)$, the training set is formed by

simply sets and not themselves random variables, this is formalized as follows:

Let F be a random variable taking values in the n -fold cartesian product space of simplices S_r . Let f be a particular instantiation of that variable, i.e., an element in the n -fold cartesian product space of simplices S_r . Then f can be viewed as a Euclidean vector, with indices given by a value $x \in \mathbf{X}$ and $y \in \mathbf{Y}$. Accordingly, we can indicate a component of f by writing $f(x, y)$. So for all x, y , $f(x, y) \geq 0$, and for any fixed x , $\sum_y f(x, y) = 1$.

This defines the random variable F . The formal sense in which this F can be viewed as an “ \mathbf{X} -conditioned distribution over \mathbf{Y} ” arises in how it is statistically related to certain other random variables (specified below) taking values in \mathbf{X} and in \mathbf{Y} .

- In a similar fashion, the generalizer’s hypothesis is an “ \mathbf{X} -conditioned distribution over \mathbf{Y} ”, i.e., the hypothesis random variable H takes values in the n -fold cartesian product space of simplices S , and components of any instantiation h of H can be indicated by $h(x, y)$.
- If for all x $h(x, y)$ is a Kronecker delta function (over y), h is called “single-valued”, and similarly for f . In such a case, the distribution in question reduces to a single-valued function from \mathbf{X} to \mathbf{Y} .
- The value d of the training set random variable is an ordered set of m input-output pairs, or “examples”. Those pairs are indicated by $d_{\mathbf{X}}(i)$, $d_{\mathbf{Y}}(i)$ $\{i = 1 \dots m\}$. The set of all input values in d is $d_{\mathbf{X}}$ and similarly for $d_{\mathbf{Y}}$. m' is the number of distinct values in $d_{\mathbf{X}}$.
- The cost C is a real-valued random variable.
- The primary random variables are such target distributions F , such hypothesis distributions H , training sets D , and real-world “cost” or “error” values C measuring how well one’s learning algorithm performs. They are “coupled” to supervised learning by imposing certain conditions on the relationship between them, conditions that are discussed next.

THE RELATIONSHIP BETWEEN C , F AND H , MEDIATED BY Q , Y_F AND Y_H

It will be useful to relate C to F and H using three other random variables. “Testing” (involved in determining the value of C) is done at the \mathbf{X} value given by the \mathbf{X} -valued random variable Q . \mathbf{Y} values associated with the hypothesis and Q are given by the \mathbf{Y} -valued random variable Y_H (with instantiations y_H), and \mathbf{Y} values associated with the target and Q are given by the \mathbf{Y} -valued random variable Y_F (with instantiations

applicable results as possible, care is taken in this appendix to discuss how a number of different learning scenarios can be cast in terms of the EBF.

NOTATION

- In general, unless indicated otherwise, random variables are written using upper case letters. A particular instantiation value of such a random variable is indicated using the corresponding lower case letter. Note though that some quantities (e.g., parameters like the size of the spaces) are neither random variables nor instantiations of random variables, so their written case carries no significance.
- When clarity is needed, the argument of a $P(\cdot)$ will not be used to indicate what the distribution is; rather a subscript will denote the distribution. For example, $P_F(h)$ means the prior over the random variable F (targets), evaluated at the value h (a particular hypothesis). This is common statistics notation. (Note that with conditioning bars, this notation leads to expressions like “ $P_{A|B}(c | d)$ ”, meaning the probability of random variable A conditioned on variable B , evaluated at values c and d , respectively.)
- Also in accord with common statistics notation, “ $E(A | b)$ ” will be used to mean the expectation value of A given $B = b$, i.e., to mean $\int da \, a P(a | b)$. (Sums replace integrals if appropriate.) This means in particular that anything not specified is averaged over. So for example, $E(A | b) = \int dc \, da \, a P(a | b, c) P(c | b) = \int dc \, E(a | b, c) P(c | b)$. When it is obvious that their value is assumed fixed and what it is fixed to, sometimes I won’t specify variables in conditioning arguments.
- I will use n and r to indicate the (countable though perhaps infinite) number of elements in the set \mathbf{X} (the input space) and the set \mathbf{Y} (the output space) respectively. (\mathbf{X} and \mathbf{Y} are the only case in this paper where capital letters don’t indicate random variables.) Such cases of countable \mathbf{X} and \mathbf{Y} are the simplest to present, and always obtain in the real world where data is measured with finite precision instruments and is manipulated on finite size digital computers.

A generic \mathbf{X} value is indicated by x , and a generic \mathbf{Y} value by y . Sometimes I will implicitly take \mathbf{Y} and/or \mathbf{X} to be sets of real numbers. (This is the case when talking about the “expected value” of a \mathbf{Y} -valued random variable for example.)

THE PRIMARY RANDOM VARIABLES

- In this paper, the “true” or “target” relationship between (test set) inputs and (test set) outputs is taken to be an \mathbf{X} -conditioned distribution over \mathbf{Y} (i.e., intuitively speaking, a $P(y | x)$). In other words, where S_r is defined as the r -dimensional unit simplex, the “target distribution” is a mapping $\mathbf{X} \rightarrow S_r$. Since \mathbf{X} and \mathbf{Y} are

$\delta(y_F, \phi(q))$. However assume that d is created by corrupting ϕ with both noise in \mathbf{X} and noise in \mathbf{Y} . This can be viewed as a “function + noise” scenario where the noise present in generating the training set is absent in testing. (This case is discussed in some detail below.)

As an example of such a scenario, viewing any particular pair of \mathbf{X} and \mathbf{Y} values from the training set as random variables X_t and Y_t , one might have $Y_t = \sum_{X'} \gamma(X_t, X') \phi(X') + \epsilon$, where X' is a dummy X variable, $\gamma(., .)$ is a convolutional process giving noise in \mathbf{X} , and ϵ is a noise process in \mathbf{Y} . (Strictly speaking, this particular kind of \mathbf{Y} -noise requires that $r = \infty$, as otherwise $\sum_{X'} \gamma(x, x') \phi(x') + \epsilon$ might not lie in \mathbf{Y} .)

For this scenario, $t \neq f$. In addition, $P(d | f)$ doesn't have the form given in equation (2.1). In particular, due to the convolution term, $P(d | f)$ will depend on the values of $f = \phi$ for $x \notin d_{\mathbf{X}}$; the likelihood for this scenario is not vertical.

6. If there are only two possible $L(., .)$ values (for example), $P(c' | d)$ does give $P(c | d)$. This is because $P(c' | d)$ gives $E(C' | d) = E(C | d)$ (see below in appendix A), and since there are two possible costs, $E(C | d)$ gives $P(c | d)$. It is for more than two possible cost values that the distributions $P(c' | d)$ and $P(c | d)$ do not determine one another. In fact, even if there are only two possible values of $L(., .)$, so that $P(c' | d)$ sets $P(c | d)$, it does not follow that $P(c | d)$ sets $P(c' | d)$. As an example, consider this case where $n = m' + 2$, and we have zero-one loss. Assume that given some d , $P(f | d)$ and $P(h | d)$ are such that either h agrees exactly with f for OTS q or the two never agree, with equal probability. This means that for zero-one OTS error, $P(c | d) = \delta(c, 0) / 2 + \delta(c, 1) / 2$. However we would get the same distribution if all four possible agreement relationships between h and f for the off-training set q were possible, with equal probabilities. And in that second case, we would have the possibility of C' values that are impossible in the first case (e.g., $c = 1/2$). QED.

7. In general, whether “error” means C or C' , this quantity is of interest only if the number of values error can have is large. So for example, it is of interest for C' if r is large and we have quadratic loss.

8. Noise in \mathbf{X} of the form mentioned in footnote 5 will not be considered in this paper. The extension to analyze such noise processes is fairly straight-forward however.

APPENDIX A. DETAILED EXPOSITION OF THE EBF

This appendix discusses the EBF in some detail. Since it is the goal of this paper to present as broadly

$\pi(x)$, many distributions behave very differently for OTS rather than IID error. (See section (2).)

Acknowledgments: I would like to thank Cullen Schaffer, Wray Buntine, Manny Knill, Tal Grossman, Bob Holte, Tom Dietterich, Karl Pfleger, Mark Plutowski, Bill Macready, Bruce Mills, and Jeff Jackson for interesting discussions. This work was supported in part by the Santa Fe Institute and by TXN Inc.

FOOTNOTES

1. For more than two possible values of $L(., .)$, it's not clear what happens. Nor is it clear how much of this carries over to costs C' (see section (2)) rather than C .

2. All of this is a formal statement of a rather profound (if somewhat philosophical) paradox: How is it that we perform inference so well in practice, given the NFL theorems and the limited scope of our prior knowledge? A discussion of some "head-to-head minimax" ideas that touch on this paradox is presented in paper two.

3. All that is being argued in this discussion of classes (1) and (2) is that the absence of a punt signal does not provide a reason to believe error is good. This argument doesn't directly address whether the presence of a punt signal gives you reason to believe you're in class (1), and therefore is correlated with bad error. The explanation of why there is no such correlation is more subtle than simply counting the number of f 's in each class. It involves the fact that there are actually a continuum of classes, and that for fixed f , raising s (so as to get a punt signal) lowers OTS (!) error.

4. In general, π itself could be a random variable that can be estimated from the data, that is perhaps coupled to other random variables (e.g., f), etc. However here I make the usual assumption in the neural net and computational learning literature that π is fixed. This is technically known as a "filter likelihood", and has powerful implications. See [Wolpert 1994b].

5. As an example, assume that f is some single-valued function from \mathbf{X} to \mathbf{Y} , ϕ , so that $P(y_F | f_\phi, q) =$

in the real world. Rather the issue is what can be formally established about how well they work in the real world without making any assumptions concerning targets.

vii) Differences between the NFL theorems and computational learning theory

Despite the foregoing, there are some similarities between the NFL theorems and computational learning theory. In particular, when all targets are allowed - as in the NFL theorems - PAC bounds on the error associated with $s = 0$ are extremely poor [Blumer et al 1987, 1989, Dietterich 1990, Wolpert 1994a]. However there are important differences between the NFL theorems and this weak-PAC-bounds phenomenon.

1) For the most part, PAC is designed to give positive results. In particular, this is the case with the PAC bounds mentioned above. (More formally, the bounds in question give an upper bound on the probability that error exceeds some value, not a lower bound.) However lack of a positive result is not the same as a negative result, and the NFL theorems are full-blown negative results.

2) PAC (and indeed all of computational learning theory) has nothing to say about this-data (i.e., Bayesian) scenarios. They only concern data-averaged quantities. PAC also is primarily concerned with polynomial versus exponential convergence issues, i.e., asymptotics of various sorts. The NFL theorems hold even if one doesn't go to the limit, and hold even for this-data scenarios. (See also [Wolpert 1994a] for a discussion of how PAC's being exclusively concerned with convergence issues renders its real-world meaningfulness debatable, at best.)

3) The PAC bounds in question can be viewed as saying there is no universally good learning algorithm. They say nothing about the possibility of whether some algorithm 1 may be better than some other algorithm 2 in most scenarios. As a particular example, nothing in the PAC literature suggests that there are as many (appropriately weighted) f 's for which a boosted learning algorithm [Drucker et al. 1993, Shapire 1990] performs *worse* than its unboosted version as there are for which the reverse is true.

4) The PAC bounds in question don't emphasize the importance of a vertical likelihood; they don't emphasize the importance of homogenous noise when the target is a single-valued function; they don't emphasize the importance of whether the loss function is homogenous; they don't invoke "scrambling" for non-homogenous loss functions (indeed, they rarely consider such loss functions); they don't concern averaging over pairs of h 's (in the sense of section (4)), etc. In all this, they are too general. Note that this over-generality extends beyond the obvious problem that they are "(sampling) distribution free". Rather they are too general in that they are independent of many of the features of a supervised learning problem that are crucially important.

5) Computational learning theory does not address OTS error. Especially when m is not infinitesimal in comparison to n and/or $\pi(x)$ is highly non-uniform, computational learning theory results are changed significantly if one uses OTS error (see [Wolpert 1994a]). And even for infinitesimal m and fairly uniform

and (2) carefully: In general, either

1) The target is such that the algorithm will almost always punt, *but when it does not punt, it usually makes significant errors*, or

2) The target is such that the algorithm has tiny expected error when it chooses not to punt.

Now there are many more f 's in class 1 than in class 2. So even though the probability of our no-punt signal is small for each of the f 's in class 1 individually, when you multiply by the number of such f , you see that the probability of being in class 1, given that you have a no-punt signal, isn't worse than the probability of being in class 2, given the same signal. In this sense, the signal gains you nothing in determining in which class you're in, and therefore in determining likely error.³

So at a minimum, one must assume that $P(f)$ is not uniform to have justification for believing the punt / no punt signal. Now one could argue that a uniform $P(f)$ is highly unlikely when there's a no-punt signal, i.e., $P(\text{no punt} \mid \alpha = \text{uniform } P(f), m)$ is very small, and that this allows one to dismiss this value of α if we see a no punt signal. Formally though, α is a hyperparameter, and should be marginalized: it is axiomatically true that $P(f) = \int d\alpha P(f \mid \alpha) P(\alpha)$ and is fixed beforehand, independent of the data. So the presence / absence of a punt signal can't be used to "infer" something about $P(f)$, formally speaking. (See the discussions of hierarchical Bayesian analysis and empirical Bayes in [Berger 1985] and [Bernardo and Smith 1994].) More generally, the NFL theorems allow us to "jump a level", so that classes 1 and 2 refer to α 's rather than f 's. And at this new level, we again run into the fact that there are many more elements in class 1 than in class 2.

To take another perspective, although the likelihood $P(\text{no punt} \mid \text{class}, m)$ strongly favors class 2, the posterior need not. Lack of appreciation for this distinction is an example of how computational learning theory relies almost exclusively on likelihood-driven calculations, ignoring posterior calculations.

It may be useful to directly contrast the intuition behind the class 1-2 reasoning and that behind the NFL theorems: The class 1-2 logic says that given an f with a non-negligible percentages of 1's, it's hugely unlikely to get all 0's in a large random data set. Hence, so this intuitive reasoning goes, if you get all 0's, you can conclude that f does not have a non-negligible percentages of 1's, and therefore you are safe in guessing 0's outside the training set. The contrasting intuition: say you are given some particular training set, say of the first K points in \mathbf{X} , together with associated \mathbf{Y} values. Say the \mathbf{Y} values happen to be all 0's. Obviously, without some assumption concerning the coupling of f 's behavior over the first K points in \mathbf{X} with its behavior outside of those points, f could have any conceivable behavior outside of those points. So the fact that it is all 0's has no significance, and cannot help you in guessing.

It should be emphasized that none of the reasoning of this subsection directly addresses the issue of whether the punting algorithm has good "head-to-head minimax" OTS behavior in some sense (see paper two). That is an issue that has yet to be thoroughly investigated. In addition, recall that no claims are being made in this paper about what is (not) reasonable in practice; punting algorithms might very well work well

have the following result, proven in appendix E:

Theorem 10 For the h^* learning algorithm, for all targets ϕ such that $\phi(x) = 0$ for more than m distinct x , $E(C_{OTS} \mid \phi, \text{punt}, m) \leq E(C_{OTS} \mid \phi, \text{no punt}, m)$.

For $n \gg m$, essentially all ϕ meet the requirement given in theorem (10); for such n and m , we do better to follow the algorithm's guessing advice when we are told not to than we are told the guessing is good!

In many respects, the proper way to analyze punting algorithms is given by decision theory. First, assign a cost to punting. (Formally, this just amounts to modifying the form of $P(c \mid f, h, d)$ for the case where h and d lead to a punt signal.) This cost should not be less than the minimal no-punting cost, or the optimal algorithm is to never guess. Similarly, it should not be more than the maximal no-punting cost, or the optimal algorithm never punts. Given such a punting cost, the analysis of a particular punting algorithm consists of finding those $P(f)$ such that $E(C_{OTS} \mid m)$ is "good" (however defined). In lieu of such an analysis, one can find those $P(f)$ such that $E(C_{OTS} \mid \text{no punt}, m) < E(C_{OTS} \mid \text{punt}, m)$. (E.g., one can analyze whether priors that are uniform in some sphere centered on h^* and zero outside of it result in this inequality.) Such analyses - apparently never carried out by proponents of punting algorithms - are beyond the scope of this paper however. (In addition, they vary from punting algorithm to punting algorithm.)

vi) Intuitive arguments concerning the NFL theorems and punting algorithms

Consider again the algorithm addressed in theorem (10). For this algorithm, there are two separate kinds of f :

- 1) f such that the algorithm will almost always punt for a d sampled from f , or
- 2) f such that the algorithm has tiny expected error when it chooses not to punt.

(Targets f with almost no x 's such that $f(x) = 1$ are in the second class, and other targets are in the first class.)

It might seem that this breakdown justifies use of the algorithm. After all, if the target is such that there's a non-negligible probability that the algorithm doesn't punt, it's not in class 1, so error will be tiny. Thus it would seem that *whatever the target* (or prior over targets), if the algorithm hasn't punted, we can be confident in its guess. (Similar arguments can be made when the two classes distinguish sets of $P(f)$'s rather than f 's.)

However if we restate this, the claim is that $E(C \mid \text{no punt}, m)$ is tiny for sufficiently large m , for any prior over targets $P(\phi)$. (Note that for $n \gg m$, m' is unlikely to be much less than m .) This would imply in particular that it is tiny for uniform $P(\phi)$. However from the preceding subsection we know that this isn't true. So we appear to have a paradox.

To resolve this paradox, consider using our algorithm and observing the no punt signal. Now restate (1)

would result in large error. To analyze this, I will slightly modify the definition of punting algorithms so that they always guess, but also always output a punt / no punt signal (and perhaps ask for more training set elements), based only on the d at hand. The issue at hand then is how the punt / no punt signal is statistically correlated with C .

Examine *any* training set d for which some particular algorithm outputs a no punt signal. By the NFL theorems, for such a d , for uniform $P(f)$, a vertical $P(d | f)$, and a homogenous OTS error, $P(c | d)$ is the same as that of a random generalizer, i.e., under those conditions, $P(c | d, \text{no punt}) = \Lambda(c) / r$. As a direct corollary, $P(c | m, \text{no punt}) = \Lambda(c) / r$. It follows that $P(c | \text{no punt}) = \Lambda(c) / r$ (assuming the no punt signal arises while OTS error is still meaningful, so $m' < n$).

Using the same kind of reasoning though, we also get $P(c | \text{punt}) = \Lambda(c) / r$, etc. So there is no statistical correlation between the value of the punt signal and OTS error. Unless we assume a non-uniform $P(f)$, even if our algorithm “grows” d until there is a no punt signal, the value of the punt / no punt signal tells us nothing about C . Similar conclusions follow from comparing a punting algorithm to its “scrambled” version, as in the analysis of non-homogenous error (see paper two).

In addition, let A and B be two punting algorithms that are identical in when they decide to output a punt signal, but B guesses randomly for all test inputs $q \notin d_X$. Then for the usual reasons, A ’s distribution over OTS error is, on average, the same as that of B , i.e., no better than random. This is true even if we condition on having a no punt signal.

One nice characteristic of some punting algorithms - the characteristic exploited by those who advocate such algorithms - is that there can be some prior-free assurances associated with them. As an example, for all targets f , the probability of such an algorithm guessing and making an error in doing so is very small (see classes (1) and (2) below): $\forall f$, for sufficiently large m and non-negligible ϵ , $P(c_{OTS} > \epsilon, \text{no punt} | f, m)$ is tiny.

However $P(c_{OTS} > \epsilon, \text{no punt} | f, m)$ in fact equals 0 for the always-punt algorithm. So one might want to also consider other distributions like $P(c_{OTS} > \epsilon | \text{no punt}, f, m)$ or $P(c_{OTS} < 1 - \epsilon, \text{no punt} | f, m)$ to get a more definitive assessment of the algorithm’s utility. Unfortunately though, both of these distributions are highly f -dependent. (This illustrates that the f -independent aspects of the punting algorithm mentioned in the previous paragraph do not give a full picture of the algorithm’s utility.)

In addition, other f -independent results hardly inspire confidence in the idea of making a guess only when there is a no punt signal. As an illustration, restrict things so that both hypotheses h and targets are single-valued (and therefore targets are written as functions ϕ), and there is no noise. \mathbf{Y} is binary, and we have zero-one loss. Let the learning algorithm always guess the all 0’s function, h^* . The punt signal is given if d_Y contains at least one non-zero value. Then for the likelihood of (2.1), uniform $\pi(x)$, and $n \gg m$, we

to them. In the real world, it is almost always the case that we know d and h in full, not simply functions of them like the empirical error. In such a scenario, it is hard to see why one would be concerned with a distribution of the form $P(c \mid \text{function}(d), h)$, as opposed to distributions of the form $P(c \mid d)$ (or perhaps $P(c \mid d, h)$, or the f -average of $P(c \mid d, f)$, or some such). So since the NFL theorems say there is no *a priori* distinction between algorithms as far as $P(c \mid d)$ is concerned, it is hard to see why one should choose between algorithms based on distributions of the form $P(c \mid \text{function}(d), h)$, if one does indeed know d in full.

iv) Implications of the NFL theorems for active learning algorithms

Active learning (aka “query-based learning”, or “membership queries”) is where the learner decides what the points $d_{\mathbf{X}}$ will be. Usually this is done dynamically; as one gets more and more training examples, one uses those examples to determine the “optimal” next choices of $d_{\mathbf{X}}(i)$.

As far as the EBF is concerned, the only difference between active learning and traditional supervised learning is in the likelihood. Rather than IID likelihoods like that in equation (2.1), in active learning each successive $d_{\mathbf{X}}(i)$ is a function of the $(i - 1)$ pairs $\{d_{\mathbf{X}}(j = 1, i - 1), d_{\mathbf{Y}}(j = 1, i - 1)\}$, with the precise functional dependence determined by the precise active learning algorithm being used.

So long as it is true that $P(d_{\mathbf{Y}}(m) \mid d_{\mathbf{X}}(m), f)$ is independent of $f(x \neq d_{\mathbf{X}}(m))$, active learning has a vertical likelihood (see appendix C). So all of the negative implications of the NFL theorems apply just as well to active learning as IID likelihood learning, and in particular apply to the kinds of active learning discussed in the computational learning community.

v) Implications of the NFL theorems for “punting” learning algorithms

Some have advocated using algorithms that have an extra option besides making a guess. This option is to “punt”, i.e., refuse to make a guess. As an example, an algorithm might choose to punt because it has low confidence in its guess (say for VC theory reasons). It might appear that, properly constructed, such algorithms could avoid making bad guesses. If this were the case, it would be an assumption-free way of ensuring that *when one guesses*, the guesses are good. (One would have traded in the ability to always make a guess to ensure that the guesses one does make are good ones.) In particular, some have advocated using algorithms that add elements to d adaptively until (and if) they can make what they consider to be a safe guess.

However the simple fact that a particular punting algorithm has a small probability of making a poor guess, by itself, is no reason to use that algorithm. After all, the completely useless algorithm that always punts has zero probability of making a poor guess. Rather what is of interest is how well the algorithm performs when it does guess, and/or how accurate its punt-signal warning is as an indicator that to make a guess

the value of s tells you nothing whatsoever about the value of c'_{OTS} . (See [Wolpert 1994a] for further discussion of how independence of s and c'_{OTS} is compatible with the VC theorems.)

Intuitively, many of the computational learning theory results relating empirical error s and generalization error c'_{IID} are driven by the fact that s is formed by sampling c'_{IID} (see [Wolpert 1994a]). However for OTS c' the empirical error s cannot be viewed as a sample of c' . Rather s and c'_{OTS} are on an equal footing. Indeed, for single-valued targets and hypotheses, and no noise, s and c'_{OTS} are both simply the value c'_{IID} has when restricted to a particular region in \mathbf{X} . (The region is $d_{\mathbf{X}}$ for s , $\mathbf{X} - d_{\mathbf{X}}$ for c'_{OTS} .) In this sense, there is symmetry between s and c'_{OTS} (symmetry absent for s and c'_{IID}). Given this, it should not be surprising that for uniform $P(f)$, the value of s tells us nothing about the value of c'_{OTS} and vice-versa.

iii) Implications for Vapnik-Chervonenkis results

The s -independence of the results presented above has strong implications for the uniform convergence formalism for investigating supervised learning [Vapnik 1982, Vapnik and Bottou 1993, Anthony and Biggs 1992, Natarajan 1991, Wolpert 1994a]. Consider zero-one loss, where the empirical error s is very low and the training set size m is very large. Assume that our learning algorithm has a very low VC dimension. Since s is low and m large, we might hope that that low VC dimension confers some assurance that our generalization error will be low, independent of assumptions concerning the target. (This is one common way people try to interpret the VC theorems.)

However according to the results presented above, low s , large m , and low VC dimension, by themselves, provide no such assurances concerning OTS error (unless one can somehow *a priori* rule out a uniform $P(f)$ - not to mention rule out any prior having even more dire implications for generalization performance). This is emphasized by the example given above where a tight confidence interval on the probability of c'_{OTS} differing from s arises solely from $P(c'_{OTS} | m)$ and $P(s | m)$ being peaked about the same value; s and c'_{OTS} are statistically independent, so knowing s tells you nothing concerning c'_{OTS} . Indeed, presuming c'_{OTS} is small due only to the fact that s , m and the learning algorithm's VC dimension are small can have disastrous real-world consequences (see the example concerning "We-Learn-It Inc." in [Wolpert 1994a]).

Of course, there are many other conditioning events one could consider besides the ones considered in this paper. And in particular, there are many such events that involve empirical errors. For example, one might investigate the behavior of the uniform f -average of $P(c | s_A, s_B, m, f)$, where s_A and s_B are the empirical errors for the two algorithms A and B considered in example (i) in section (2).

It may well be that for some of these alternative conditioning events involving empirical errors, one can find *a priori* distinctions between learning algorithms, dependences on s values, or the like. Although such results would certainly be interesting, one should be careful not to ascribe too much practical significance

raised, this does not mean that the integral exhibits the same behavior.

As an aside, it should be noted that the only property of s needed by theorem (9) or its corollary is that $P(s | d, f) = P(s | d)$. In addition to holding for the random variable S , this property will hold for any random variable σ that is a function only of d for the algorithms under consideration. So in particular, we can consider using cross-validation to choose among a set of one or more deterministic algorithms. Define σ as the cross-validation errors of the algorithms involved. Since for a fixed set of deterministic algorithms σ is a function solely of d , we see that for a uniform $P(f)$, σ is statistically independent from C ; there is no information contained in the set of cross-validation errors that has bearing on generalization error. In this sense, unless one makes an explicit assumption for $P(f)$, cross-validation error has no use as an estimate of generalization error.

ii) Compatibility with Vapnik-Chervonenkis results

The fact that $P(c | s, m) = \Lambda(c) / r$ (under the appropriate conditions) means that $P(c | s, m) = P(c | m)$ under those conditions (see corollary (1)). This implies that $P(s | c, m)$ is independent of cost c . So C_{OTS} and empirical error S are statistically independent, for uniform $P(f)$, homogenous L , and a vertical likelihood. Indeed, in appendix B in [Wolpert 1992] there is an analysis of the case where we have a uniform sampling distribution $\pi(\cdot)$, zero-one loss, binary \mathbf{Y} , and $n/m \rightarrow \infty$ (so $C''_{OTS} \rightarrow C'_{IID}$; see appendix B of this paper). It is there proven that $E(C'_{IID} | s, m) = 1/2$, independent of s .

In accord with this, one expects that C'_{OTS} and S are independent for uniform $P(f)$. On the other hand, Vapnik-Chervonenkis (uniform convergence) theory tells us that $P(c'_{IID} - s | m)$ is biased towards small values of $c'_{IID} - s$ for low-VC dimension generalizers, and large m . This is true for any prior $P(f)$, and therefore in particular for a uniform prior. It is also true even when $n \gg m$, so that C'_{OTS} and C'_{IID} closely approximate each other.

It should be emphasized that there is no contradiction between these VC results and the NFL theorems. Independence of s and c'_{OTS} does not imply that s and c'_{OTS} can differ significantly. For example, both the VC results and the NFL theorems would hold if $P(c'_{OTS} | m)$ and $P(s | m)$ were independent but were both tightly clumped around the same value, ζ .

Now let's say we have an instance such a "clumping" phenomenon, but don't know ζ (ζ being determined by (the unknown) f , among other things). We might be tempted to take the observed value of s as an indicator of the likely value of ζ . In turn, we might wish to view this likely value of ζ as an indicator of the likely value of c'_{OTS} . In this way, having observed a particular value of s , we could infer something about c'_{OTS} (e.g., that it is unlikely to differ from that observed value of s). However theorem (9) says that this reasoning is illegal (at least for uniform $P(f)$). Statistical independence is statistical independence; knowing

orem (1) (that holds even for non-deterministic learning algorithms, capable of guessing non-single-valued h 's):

For homogenous loss L , the uniform average over all f of $P(c | f, d, s)$ equals $\Lambda(c) / r$.

Unfortunately, one cannot continue paralleling the analysis in section (2) past this point, to evaluate quantities like the uniform average over all f of $P(c | f, s, m)$. The problem is that whereas $P(d | f, m)$ is independent of $f(x \notin d_X)$ (for a vertical likelihood), the same need not be true of $P(d | f, s, m)$. Indeed, often there are f 's for which $P(c_{OTS} | f, s, m)$ is not defined; for no d sampled from that f will an h be produced that has error s with that d . In such scenarios the uniform f -average of $P(c_{OTS} | f, s, m)$ is not defined. Moreover, the set of f for which $P(c_{OTS} | f, s, m)$ is defined may vary with s . The repercussions of this carry through for any attempt to create s -conditioned analogs of the NFL theorems. (An example of how the NFL theorems need not hold for s -conditioned distributions is presented in appendix C.)

In fact, it is hard to say anything general about $P(c | f, s, m)$. In particular, it is not always the case that higher s results in lower c_{OTS} if f is fixed, as in the example in appendix C. To see this, consider the scenario given there with a simple change in the learning algorithm. For the new learning algorithm, if all input elements of the training set, d_X , are in some region Ξ , then an hypothesis h is produced that happens to equal the target f , whereas for any other d_X 's, there are errors both on and off d_X . So if $s = 0$, we know that $c_{OTS} = 0$. But if $s > 0$, we know that $c_{OTS} > 0$; raising s from 0 has raised expected C_{OTS} .

Now consider $P(c | s, d)$ for uniform $P(f)$, where it is implicitly assumed that for at least one h for which $P(h | d) \neq 0$, the empirical error is s , so $P(s, d) \neq 0$. For this quantity we do have an NFL result that holds for any learning algorithm (see appendix C):

Theorem 9 For homogenous L , OTS error, a vertical likelihood, and uniform $P(f)$, $P(c | s, d) = \Lambda(c) / r$.

The immediate corollary is that for homogenous L , OTS error, a vertical likelihood, and uniform $P(f)$, $P(c | s, m) = \Lambda(c) / r$, independent of the learning algorithm.

It is interesting to note that a uniform $P(f)$ can give NFL for $P(c | s, m)$ even though a uniform average over f of $P(c | f, s, m)$ does not. This illustrates that one should exercise care in equating the basis of NFL for f -conditioned distributions (theorem (2)) with having a uniform prior.

An immediate question is how theorem (9) can hold despite the example above where as s shrinks $E(C_{OTS} | f, s, m)$ grows, for any target f . The answer is that $P(c | s, m) = \int df P(c | f, s, m) P(f | s, m)$. Even if for any *fixed* target f the quantity $P(c | f, s, m)$ gets biased towards lower cost c as the empirical error s is

4. THE NFL THEOREMS AND COMPUTATIONAL LEARNING THEORY

This section discusses the NFL theorem's implications for and relationship with computational learning theory.

Define the empirical error

$$s \equiv \frac{\sum_{i=1}^m \sum_{y_H} L(y_H, d_{\mathbf{Y}}(i)) h(d_{\mathbf{X}}(i), y_H) \pi(d_{\mathbf{X}}(i))}{\sum_{i=1}^m \pi(d_{\mathbf{X}}(i))}.$$

As an example, for zero-one loss and single-valued h , s is the average misclassification rate of h over the training set. Note that the empirical error is implicitly a function of d and h but of nothing else (π being fixed). (For deterministic learning algorithms, this reduces to being a function only of d .) So for example $P(s | d, f) = \int dh P(s | d, f, h) P(h | d) = \int dh P(s | d, h) P(h | d) = P(s | d)$.

This section first analyzes distributions over C that involve the value of s , as most of computational learning theory does. Then it analyzes OTS behavior of “membership queries” algorithms and “punting” algorithms (those that may refuse to make a guess), algorithms that are also analyzed in computational learning theory.

i) NFL theorems involving empirical error

Some of the NFL theorems carry over essentially unchanged if one conditions on s in the distribution of interest. This should not be too surprising. For example, consider the most common kind of learning algorithms, deterministic ones that produce single-valued h 's. For such learning algorithms, the training set d determines the hypothesis h and therefore determines s . So specifying s in addition to d in the conditioning statement of the probability distribution provides no information not already contained in d . This simple fact establishes the NFL theorem for $P(c | f, d, s)$, for these kinds of learning algorithms.

More generally, first follow along with the derivation of lemma (1), to get

Lemma 2 $P(c | f, d, s) = \sum_{y_H, y_F, q} \delta(c, L(y_H, y_F)) P(y_H | q, d, s) P(y_F | q, f) P(q | d),$

where use was made of the identities $P(y_F | q, f, s) = P(y_F | q, f)$, and $P(q | d, s) = P(q | d)$. (Both identities follow from the fact that $P_{A|B,S,D,H}(a | b, s(d, h), d, h) = P(a | b, d, h)$ for any variables A and B .)

Continuing along with the logic that resulted in theorem (1), we arrive at the following analogue of the-

well in the past, this provides no assurances whatsoever that you can successfully apply that knowledge to some current inference problem. The fact that a learning algorithm has been used many times with great success provides no formal (!) assurances about its behavior in the future.² After all, assuming that how well you generalized in the past carries over to the present is formally equivalent to (a variant of) cross-validation - in both cases, one tries to extrapolate from generalization accuracy on input points for which we now know what the correct answer was, to generalization behavior in general.

Finally, it is important to emphasize that results based on averaging uniformly over $f/\phi/P(\phi)$ shouldn't be viewed as normative. The uniform averaging enables us to reach conclusions that assumptions are needed to distinguish between algorithms, not that algorithms can be (profitably) distinguished without any assumptions. I.e., if such an average ends up favoring algorithm A over B (as it might for a non-homogenous loss function, for example), that only means one "should" use A if one has reason to believe that all f are equally likely *a priori*.

vii) Other peculiar properties associated with OTS error

There are many other aspects of OTS error which, although not actually NFL theorems, can nonetheless be surprising. An example is that in certain situations the expected OTS error grows as the size of the training set increases, even if one uses the best possible learning algorithm, the Bayes-optimal learning algorithm (i.e., the learning algorithm which minimizes $E(C | d)$ - see [Wolpert 1994a]). In other words, sometimes the more data you have, the less you know about the OTS behavior of ϕ , on average.

In addition, the NFL theorems have strong implications for the common use of a "test set" or "validation set" T to compare the efficacy of different learning algorithms. The conventional view is that the error measured on such a set is a sample of the full generalization error. As such, the only problem with using error on T to estimate "full error" is that error on T is subject to statistical fluctuations, fluctuations that are small if T is large enough. However if we are interested in the error for $x \notin \{d \cup T\}$, the NFL theorems tell us that (in the absence of prior assumptions) error on T is meaningless, no matter how many elements there are in T .

Moreover, as pointed out in section (3) of the second of this pair of papers, use of test sets cannot correspond to an assumption only about targets (i.e., there is no $P(f)$ that, by itself, justifies the use of test sets.) Rather use of test sets corresponds to an assumption about both targets and the algorithms the test set is being used to choose between. Use of test sets will give incorrect results unless one has a particular relationship between the target and the learning algorithms being chosen between.

In all this, even the ubiquitous use of test sets is unjustified (unless one makes assumptions). For a discussion of this point and of intuitive arguments for why the NFL theorems hold, see [Wolpert 1994a].

performs worse than pure randomness as for which it performs better. (Recall the discussion just below theorem (8).)

So for the learning scenarios considered in this section (zero-one loss, etc.) the burden is on the user of a particular learning algorithm. Unless they can somehow show that $P(\phi)$ is one of the ones for which their algorithm does better than random, rather than one of the ones for which it does worse, they cannot claim to have any formal justification for their learning algorithm.

In fact if you press them, you find that in practice very often people's assumption don't concern $P(\phi)$ at all, but rather boil down to the statement "okay, my algorithm corresponds to an assumption about the prior over targets; I make that assumption". This is unsatisfying enough a formal justification as it stands. Unfortunately though, for many algorithms, no one has even tried to write down that set of $P(\phi)$ for which their algorithm works well. This puts the purveyors of such statements in the awkward position of invoking an unknown assumption. (Moreover, for some algorithms one can show that there is *no* assumption solely concerning targets that justifies that algorithm in all contexts. This is true of cross-validation, for example; see paper two.)

Given this breadth of the implications of the uniform-average cases, it is not surprising that uniform distributions have been used before to see what one can say *a priori* about a particular learning scenario. For example, the "Ugly Duckling Theorem" [Watanabe 1985] can be viewed as (implicitly) based on a uniform distribution. Another use of a uniform distribution, more closely related to the uniform distributions occurring in this paper, appears in the "problem-averaging" work of Hughes [Hughes 1968]. (See [Waller and Jain 1978] as well for a modern view of the work of Hughes.) The words of Duda and Hart [1973] describing that work are just as appropriate here: "Of course, choosing the a priori distribution is a delicate matter. We would like to choose a distribution corresponding to the class of problems we typically encounter, but there is no obvious way to do that. A bold approach is merely to assume that problems are 'uniformly distributed'. Let us consider some of the implications (of such an assumption)..."

In this regard, note that you really would need a proof based completely on first principles to formally justify some particular (non-uniform) $P(f)$. In particular, you cannot use your "prior knowledge" (e.g., that targets tend to be smooth, that Occam's razor usually works, etc.) to set $P(f)$, without making additional assumptions about the applicability of that "knowledge" to future supervised learning problems. This is because that "prior knowledge" is ultimately an encapsulation of two things: the data set of your experiences since birth, and the data set of your genome's experiences in the several billion years it's been evolving. So if you are confronted with a situation differing at all (!) from the previous experiences of you and/or your genome, then you are in an OTS scenario. Therefore the NFL theorems apply, and you have no formal justification for presuming that your "prior knowledge" will apply off-training set (i.e., in the future).

An important example of this is the fact that even if your prior knowledge allowed you to generalize

$P(c | f, m)$ rather than averages over α of $P(c | m, \alpha)$. This is because there is no such thing as a “uniform f ” that we can restrict the average away from with the same kind of implications as restricting an average away from a uniform prior. However, by theorem (2), for any pair of algorithms, there are targets that “favor” the first of the two algorithms, and there are targets that favor of the second. So by choosing from both sets of targets, we can construct many distributions $\Gamma(f)$ that have a small support and such that the average of $P(c | f, m)$ according to $\Gamma(f)$ is the same for both algorithms. Indeed, an interesting open question is characterizing the set of such $\Gamma(f)$ for any particular pair of algorithms.

vi) On uniform averaging

The results of the preceding subsection notwithstanding, it is natural to pay a lot of attention to the original uniform average forms of the NFL theorems. When considering those forms, it should be kept in mind that the uniform averages over f (or ϕ , or $P(\phi)$) weren’t chosen because there is strong reason to believe that all f are equally likely to arise in practice. Indeed, in many respects it is absurd to ascribe such a uniformity over possible targets to the real world. Rather the uniform sums were chosen because such sums are a useful theoretical tool with which to analyze supervised learning.

For example, the implication of the NFL theorems that there is no such thing as a general-purpose learning algorithm that works optimally for all $f / P(\phi)$ is not too surprising. However even if one already believed this implication, one might still have presumed that there are algorithms that *usually* do well and those that *usually* do poorly, and that one could perhaps choose two algorithms so that the first algorithm is usually superior to the second. The NFL theorems show that this is not the case. If all f ’s are weighted by their associated probability of error, then for any two algorithms A and B there are exactly as many f ’s for which algorithm A beats algorithm B as vice-versa.

Now if one changes the weighting over f ’s to not be according to the algorithm’s probability of error, then this result would change, and one would have *a priori* distinctions between algorithms. However, *a priori*, the change in the result could just as easily favor either A or B. Accordingly, claims that “in the real world $P(f)$ is not uniform, so the NFL results don’t apply to my favorite learning algorithm” are misguided at best. Unless you can prove that the non-uniformity in $P(f)$ is well-matched to your favorite learning algorithm (rather than being “anti-matched” to it), the fact that $P(f)$ may be non-uniform, by itself, provides no justification whatsoever for your use of that learning algorithm. (See the inner product formula, theorem (1), in [Wolpert 1994a].)

In fact, the NFL theorems for averages over priors $P(\phi)$ say (loosely speaking) that there are exactly as many priors for which any learning algorithm A beats any algorithm B as vice-versa. So uniform distributions over targets are not an atypical, pathological case, out at the edge of the space. Rather they and their associated results are the average case (!). There are just as many priors for which your favorite algorithm

tion-conditioning quantity (e.g., an f) for which that algorithm is optimal (i.e., for which that algorithm beats all other algorithms), but some algorithms are not optimal for any value of such a quantity; and more generally 3) for some pairs of algorithms the NFL theorems may be met by having comparatively many targets in which algorithm A is just slightly worse than algorithm B, and comparatively few targets in which algorithm A beats algorithm B by a lot. These points are returned to in paper 2.

v) Extensions for non-uniform averaging

The uniform sums over f (or ϕ , or $P(\phi)$) in the NFL theorems are not necessary conditions for those theorems to hold. As an example, consider the version of the theorems for which targets are single-valued functions ϕ from \mathbf{X} to \mathbf{Y} , perhaps with output-space noise superimposed, and where one average over priors α . It turns out that we recover the NFL result for that scenario if we average according to any distribution over the α which is invariant under relabelling of the ϕ . We do not need to average according to the uniform distribution, and in fact can disallow all priors that are too close to the uniform prior.

More formally, we have the following variant of theorem (7), proven in appendix C:

Corollary 3 Assume OTS error, a vertical $P(d | \phi)$, homogenous loss L , and a homogenous test set noise process. Let α index the priors $P(\phi)$, and let $G(\alpha)$ be a distribution over α . Assume $G(\alpha)$ is invariant under the transformation of the priors α induced by relabelling the targets ϕ . Then the average according to $G(\alpha)$ of $P(c | m, \alpha)$ equals $\Lambda(c) / r$.

As a particular example of this result, define α^* to be the uniform prior, that is the vector all of whose components are equal. Then one $G(\alpha)$ that meets the assumption in corollary (3) is the one that is constant over α except that it excludes all vectors α lying within some L^2 distance of α^* . (I.e., one $G(\alpha)$ that meets the assumption is the one that excludes all priors α that are too close to being uniform.) This is because rearranging the components of a vector does not change the distance between that vector and α^* , so any $G(\alpha)$ that depends only on that distance obeys the assumption in corollary (3).

Combined with corollary (3), this means that $G(\alpha)$ can have structure - it can have a huge amount of structure - and we still get NFL. Alternatively, the set of allowed priors can be tiny, and restricted to priors α with a lot of structure (i.e., to priors lying far from the uniform prior), and we still get NFL. Loosely speaking, there are just as many priors that have lots of structure for which your favorite algorithm performs worse than randomly as there are for which it performs better than randomly.

An open question is whether the condition on $G(\alpha)$ in corollary (3) is a necessary condition to have the average according to $G(\alpha)$ of $P(c | m, \alpha)$ equal $\Lambda(c) / r$.

Interestingly, we do not have the same kind of result when considering averages over targets f of

spite formal arguments some have offered trying to prove the validity of some of these heuristics.

iv) General implications of the NFL theorems

The primary importance of the NFL theorems is their implication that, for any two learning algorithms A and B, according to any of the distributions $P(c | d)$, $P(c | m)$, $P(c | f, d)$ or $P(c | f, m)$, there are just as many situations (appropriately weighted) in which algorithm A is superior to algorithm B as vice-versa. So in particular, if we know that learning algorithm A is superior to B averaged over some set of targets F, then the NFL theorems tell us that B must be superior to A if one averages over all targets not in F. This is true even if algorithm B is the algorithm of purely random guessing.

Note that much of computational learning theory, much of sampling theory statistics (e.g., bias + variance results), etc., is based on quantities like $P(c | f, m)$, or on other quantities determined by $P(c | f, m)$. (See [Wolpert 1994a].) Similarly, conventional Bayesian analysis is concerned with $P(c | d)$. All of these quantities are addressed in the NFL theorems.

As a special case of the theorems, when there are only two possible values of $L(., .)$, any two algorithms are even more tightly matched in behavior than theorems (1) through (8) indicate. (An example of such an $L(., .)$ is zero-one loss, for which there are only two possible values of $L(., .)$ regardless of r .) Let C_1 and C_2 be the costs associated with two learning algorithms. Now $P(c_1 | \text{stuff}) = \sum_{c_2} P(c_1, c_2 | \text{stuff})$, and similarly for $P(c_2 | \text{stuff})$. (Examples of “stuff” are $\{d, f\}$, $\{m\}$, etc.) If $L(., .)$ can take on two values, this provides us four equations (one each for the two possible values of c_1 and the two possible values of c_2) in four unknowns ($P(c_1, c_2 | \text{stuff})$ for the four possible values of c_1 and c_2). Accordingly, if we know both $P(c_1 | \text{stuff})$ and $P(c_2 | \text{stuff})$ for both possible values of c_1 and c_2 , we can solve for $P(c_1, c_2 | \text{stuff})$ (sometimes up to some overall unspecified parameters, since our four equations are not independent). In particular, if we know that $P_{C_1}(c | \text{stuff}) = P_{C_2}(c | \text{stuff})$, then $P(c_1, c_2 | \text{stuff})$ must be a symmetric function of c_1 and c_2 . So for all of the “stuff”s in the NFL theorems, when $L(., .)$ can take on two possible values, for any two learning algorithms, $P(c_1, c_2 | \text{stuff})$ is a symmetric function of c_1 and c_2 (under the appropriate uniform average).¹

All of the foregoing applies to more than just OTS error. In general IID error can be expressed as a linear combination of OTS error plus on-training set error, where the combination coefficients depend only on $d_{\mathbf{x}}$ and $\pi(x \in d_{\mathbf{x}})$. So generically, if two algorithms have the same on-training set behavior (e.g., they reproduce d exactly), the NFL theorems apply to their IID errors as well as their OTS set errors. (See also appendix B.)

Notwithstanding the NFL theorems though, learning algorithms can differ in that: 1) for particular f , or particular (non-uniform) $P(f)$, different algorithms can have different probabilities of error (this is why some algorithms tend to perform better than others in the real world); 2) for some algorithms there is a distribu-

.7. So averaged over all those ϕ 's, we get $[1(0) + 5(.2) + 10(.7)] / [1 + 5 + 10] = .5$. This is exactly the same expected error as algorithm B has: expected error for B is $[1(1) + 5(.8) + 10(.3)] / 16 = .5$. QED.

See example 5 in paper two for a related example.

Example 2) An algorithm that uses cross-validation to choose amongst a pre-fixed set of learning algorithms does no better on average than one that doesn't, so long as the loss function is homogenous. In addition, cross-validation does no better than anti-cross-validation (choosing the learning algorithm with the *worst* cross-validation error) on average. In particular, the error on the validation set can be measured using a non-homogenous loss (e.g., quadratic loss), and this result will still hold; all that is required is that we use a homogenous loss to measure error on the test set.

Alternatively, construct the following algorithm: "If cross-validation says one of the algorithms under consideration has particularly low error in comparison to the other, use that algorithm. Otherwise, choose randomly among the algorithms." Averaged over all targets, this algorithm will do exactly as well as the algorithm that always guesses randomly among the algorithms. In this particular sense, you cannot rely on cross-validation's error estimate (unless you impose a prior over targets or some such).

Note that these results don't directly address the issue of how accurate cross-validation is as an estimator of generalization accuracy; the object of concern here is instead the error that accompanies use of cross-validation. For a recent discussion of the accuracy question (though in a non-OTS context), see [Plutowski et al 1994]. For a more general discussion of how error and accuracy-as-an-estimator are statistically related (especially when that accuracy is expressed as a confidence interval), see [Wolpert 1994a]. The issue of how accurate cross-validation is as an estimator of generalization accuracy is also addressed in the discussion just below theorem (9) below, and in the fixed-f results in paper 2.

Example 3) Assume you're a Bayesian, and calculate the Bayes-optimal guess assuming a particular $P(f)$. (I.e., you use the $P(h | d)$ that would minimize the data-conditioned risk $E(C | d)$, if your assumed $P(f)$ were the correct $P(f)$.) You now compare your guess to that made by someone who uses a non-Bayesian method. Then the NFL theorems mean (loosely speaking) that there are as many actual priors (your assumed prior being fixed) in which the other person has a lower data-conditioned risk as there are for which your risk is lower.

Example 4) Consider any of the heuristics that people have come up with for supervised learning: avoid "over-fitting", prefer "simpler" to more "complex" models, "boost" your algorithm, "bag" it, etc. The NFL theorems say that all such heuristics fail as often (appropriately weighted) as they succeed. This is true de-

so for any other d and ϕ , $P(c | \phi, d)$ is meaningless). In such a situation, averaging over all ϕ 's with d fixed (as in theorem (4)) is not well-defined. Such situations can, at the expense of extra work, be dealt with explicitly. (The result is essentially that all of the results of this section save theorem (4) are obtained.) Alternatively, one can usually approximate the analysis for such noise processes arbitrarily well by using other, infinitesimally different noise processes, processes for which $P(c | \phi, d)$ is always defined.

iii) Examples

Example 1: Say we have no noise, and the zero-one loss $L(., .)$. Fix two possible (single-valued) hypotheses, h_1 and h_2 . Let learning algorithm A take in the training set d , and guess whichever of h_1 and h_2 agrees with d more often (the “majority” algorithm). Let algorithm B guess whichever of h_1 and h_2 agrees *less* often with d (the “anti-majority” algorithm). If h_1 and h_2 agree equally often with d , both algorithms choose randomly between them. Then averaged over all target functions ϕ , $E(C | \phi, m)$ is the same for A and B.

As an example, take $n = 5$ and $r = 2$ (i.e., $\mathbf{X} = \{0, 1, 2, 3, 4\}$, and $\mathbf{Y} = \{0, 1\}$) and a uniform sampling distribution $\pi(x)$. Take m' , the number of distinct elements in the training set, to equal 4. For expository purposes, I will explicitly show that the average over all ϕ of $E(C | \phi, m')$ is the same for A and B. (To calculate the average over all ϕ of $E(C | \phi, m)$, one sums the average of $E(C | \phi, m')$ $P(m' | m)$ over all m' .) I will take h_1 = the all 1's h , and h_2 = the all 0's h .

- i) There is one ϕ that is all 0's (i.e., for which for all \mathbf{X} values, $\mathbf{Y} = 0$). For that ϕ , algorithm A always picks h_2 , and therefore $E(C | \phi, m' = 4) = 0$; algorithm A performs perfectly. For algorithm B, expected $C = 1$.
- ii) There are five ϕ 's with one 1. For each such ϕ , the probability that the training set has all four zeroes is .2. The value of C for such training sets is 1 for algorithm A, 0 for B. For all other training sets, $C = 0$ for algorithm A, and 1 for algorithm B. So for each of these ϕ 's, the expected value of C is $.2(1) + .8(0) = .2$ for A, and $.2(0) + .8(1) = .8$ for B.
- iii) There are ten ϕ 's with two 1's. For each such ϕ , there is a .4 probability that the training set has one 1, and a .6 probability that it has both 1's. (It can't have no 1's.) If the training set has a single 1, so does the OTS region, and $C = 1$ for A, 0 for B. If the training set has two 1's, then our algorithms say guess randomly, so (expected) $C = .5$ for both algorithms. Therefore for each of these ϕ 's, expected $C = .4(1) + .6(.5) = .7$ for algorithm A, and $.4(0) + .6(.5) = .3$ for B.
- iv) The case of ϕ 's with two 1's is the same as the case of ϕ 's with three 1's (just with '1' replaced by '0' throughout). Similarly, one 1 = four 1's, and zero 1's = five 1's. So it suffices to just consider the cases already investigated, where the number of 1's is zero, one, or two.
- v) Adding them up, for algorithm A we have one ϕ with (expected) $C = 0$, five with $C = .2$ and 10 with $C =$

uniform $P(\phi)$, $P(c | m)$ equals $\Lambda(c) / r$.

We are now in a position to extend the NFL theorems to the case where neither the prior nor the target is specified in the conditioning event of our distribution of interest, and the prior need not be uniform. For such a case, the NFL results concern uniformly averaging over priors $P(\phi)$ rather than over target functions ϕ .

Since there are r^n possible single-valued ϕ , $P(\phi)$ is an r^n -dimensional real-valued vector lying on the unit simplex. Indicate that vector as α , and one of its components (i.e., $P(\phi)$ for one ϕ) as α_ϕ . (More formally, α is a hyperparameter: $P(\phi | \alpha) \equiv \alpha_\phi$.) So the uniform average over all α of $P(c | m, \alpha)$ is (proportional to) $\int d\alpha P(c | m, \alpha) = \int d\alpha [\sum_\phi P(\phi | \alpha) P(c | m, \alpha, \phi)]$, where the integral is restricted to the r^n -dimensional simplex. (α is restricted to lie on that simplex, since $\sum_\phi P(\phi | \alpha) = \sum_\phi \alpha_\phi = 1$.) It is now straight-forward to use theorem (5) to establish the following result (see appendix C):

Theorem 7 Assume OTS error, a vertical $P(d | \phi)$, homogenous loss L , and a homogenous test set noise process. Let α index the priors $P(\phi)$. Then the uniform average over all α of $P(c | m, \alpha)$ equals $\Lambda(c) / r$.

It is somewhat more involved to calculate the uniform average over all priors (indexed by) α of $P(c | d, \alpha)$. The result is derived in appendix D:

Theorem 8 Assume OTS error, a vertical $P(d | \phi)$, homogenous loss L , and a homogenous test set noise process. Let α index the priors $P(\phi)$. Then the uniform average over all α of $P(c | d, \alpha)$ equals $\Lambda(c) / r$.

By corollary (2), theorem (7) means that the average over all priors α of $P(c | m, \alpha)$ equals $P(c | m, \text{uniform prior})$. Similarly, by theorem (6), theorem (8) means that the average over all priors α of $P(c | d, \alpha)$ equals $P(c | d, \text{uniform prior})$. In this sense, whatever one's learning algorithm, one can just as readily have a prior which gives worse performance than that associated with the uniform prior as one which gives better performance.

To put this even more strongly, consider again the uniform-random learning algorithm discussed at the beginning of this section. By theorems (7) and (8), for any learning algorithm, one can just as readily have a prior for which that algorithm performs worse than the random learning algorithm - *worse than random guessing* - as a prior for which one's algorithm performs better than the random learning algorithm.

It may be that for some particular (homogenous) noise process, for some training sets d and target functions ϕ , $P(c | \phi, d)$ is not defined. This is the situation for example when there is no noise (d must lie on ϕ ,

skewed and sends all $\phi(q)$ to some fixed value y_1 , then the ϕ -average is 1 for $y_F = y_1$, 0 otherwise. Intuitively, if the noise process always results in the test value y_1 , then we **can** make *a priori* distinctions between learning algorithms; an algorithm that always guesses y_1 outside of d_X will beat one that does not.

So for simplicity restrict attention to those noise processes for which the uniform ϕ -average of $P(y_F | q, \phi)$ is independent of the target output value y_F . Recall that such a (test set) noise process is called “homogenous”. If we sum our ϕ -average of $P(y_F | q, \phi)$ over all y_F , then by pulling the sum over y_F inside the average over ϕ , we see that the sum must equal 1. (Again, see the proof of the “implication of lemma (1)”.) Accordingly, the ϕ -average equals $1 / r$. So we have the following analog of theorem (1):

Theorem 4 For homogenous loss L and a homogenous test-set noise process, the uniform average over all single-valued target functions ϕ of $P(c | \phi, d)$ equals $\Lambda(c) / r$.

Note that the noise process involved in generating the training set is irrelevant to this result. (Recall that “homogenous noise” refers to y_F and y_H , and that y_F and y_H are Y values for the test process, not the training process.) This is also true for the results presented below. So in particular, all these results hold for any noise in the generation of the training set, so long as our error measure is concerned with whether or not h equals the (homogenous noise corrupted) sample of the underlying ϕ at the test point q . (Note in particular that such a measure is always appropriate for noise-free - and therefore trivially homogenous - test set generation).

We can proceed from theorem (4) to get a result for $P(c | f, m)$ in the exact same manner as theorem (1) gave theorem (2). The result is

Theorem 5 For OTS error, a vertical $P(d | \phi)$, homogenous loss L , and a homogenous test-set noise process, the uniform average over all single-valued target functions ϕ of $P(c | \phi, m)$ equals $\Lambda(c) / r$.

Just as the logic behind theorem (2) also resulted in theorem (3), so we can use the logic behind theorem (5) to derive the following.

Theorem 6 For OTS error, a vertical $P(d | \phi)$, homogenous loss L , uniform $P(\phi)$, and a homogenous test-set noise process, $P(c | d)$ equals $\Lambda(c) / r$.

Just as theorem (3) resulted in corollary (1), so theorem (6) establishes the following.

Corollary 2 For OTS error, vertical $P(d | \phi)$, homogenous loss L , a homogenous test-set noise process, and

guessing, but rather that our algorithm performs worse than randomly!

Using similar reasoning to that used to prove theorem (2), we can derive the following theorem concerning the distribution of interest in conventional Bayesian analysis, $P(c | d)$:

Theorem 3 For OTS error, a vertical $P(d | f)$, uniform $P(f)$, and a homogenous loss L , $P(c | d) = \Lambda(c) / r$.

The reader should be wary of equating the underlying logic behind a target-averaging NFL theorem (e.g., theorem (2)) with that behind a uniform-prior NFL theorem (e.g., theorem (3)). In particular, there are scenarios (i.e., conditioning events in “ $P(c | \dots)$ ”) in which one of these kinds of NFL theorem holds but not the other. See the discussion surrounding theorem (9) below for an example.

As an immediate corollary of theorem (3), we have the following.

Corollary 1. For OTS error, a vertical $P(d | f)$, uniform $P(f)$, and a homogenous loss L , $P(c | m) = \Lambda(c) / r$.

As an aside, so long as $L(a, b) = L(b, a)$ for all pairs a and b , the mathematics of the EBF is symmetric under interchange of h and f . (In particular, for any loss L , $P(f | h, d) = P(f | d)$, and $P(h | f, d) = P(h | d)$.) Accordingly, all of the NFL theorems have analogues where the hypothesis h rather than the target f is fixed and then uniformly averaged over. So for example, for OTS error, homogenous $L(., .)$, and a generalizer such that $P(d | h)$ is independent of $h(x \notin d_X)$, the uniform average over h of $P(c | h, m) = \Lambda(c) / r$. (For such a generalizer, assuming $h_1(x) = h_2(x)$ for all $x \in d_X$, the probability that the training set used to produce the hypothesis was d is the same, whether that produced hypothesis is h_1 or h_2 .) Such results say that averaged over all h ’s the algorithm might produce, all posteriors over targets (and therefore all priors) lead to the same probability of cost, under the specified conditions.

ii) Averaging over all functions ϕ

Now consider the scenario where only those targets f are allowed that can be viewed as single-valued functions ϕ from \mathbf{X} to \mathbf{Y} with noise superimposed (see section 2). To analyze such a scenario, I will no longer consider uniform averages involving f directly, but rather uniform averages involving ϕ . Accordingly, such averages are now sums rather than integrals. (For reasons of space, only here in this subsection will I explicitly consider the case of f ’s that are single-valued ϕ ’s with noise superimposed.)

In this new scenario, lemma (1) still holds, with f replaced by ϕ . However now we cannot simply set the uniform ϕ -average of $P(y_F | q, \phi)$ to $1 / r$, in analogy to the reasoning used above (see the proof of the “implication of lemma (1)” in appendix C). To give an extreme example, if the test set noise process is highly

As an example of theorem (1), for zero-one loss, we get the f-average of $E(C | f, d) = 1 / r$. More generally, for an even broader set of loss functions L than homogenous L 's, the sum over target outputs y_F of $L(y_H, y_F)$ is independent of the hypothesis output, y_H . For such L 's we get generalizer-independence for the uniform average over targets f of $E(C | f, d)$, even if we don't have such independence for the uniform average of $P(c | f, d)$.

Note that theorem (1) doesn't rely on having q lie outside of d_X ; it holds even for IID error. In addition, since both f and d are fixed in the conditional probability in theorem (1), any statistical coupling between f and d is ignored in that theorem. For these kinds of reasons, theorem (1) isn't too interesting by itself. The main use of it is to derive other results, results that rely on using OTS error and that are affected by the coupling of targets f and training sets d . As the first of these, I will show how to use theorem (1) to evaluate the uniform f-average of $P(c | f, m)$ for OTS error.

In evaluating the uniform f-average of $P(c | f, m)$, not all f 's contribute the same amount to the answer. That is because

$$P(c | f, m) = \sum_d P(c | f, d) P(d | f),$$

and so long as the likelihood $P(d | f)$ is not uniform over f , we cannot just pull the outside f-average through to reduce the $P(c | f, d)$ to $\Lambda(c) / r$. This might lead one to suspect that if the learning algorithm is "biased" towards the targets f contributing the most to the uniform f-average of $P(c | f, m)$, then the average would be weighted towards low values of cost, c . However this is wrong; it turns out that the uniform f-average of $P(c | f, m)$ is independent of the learning algorithm, if one restricts oneself to OTS error.

In fact, assume that we have any $P(q | d)$ such that $P(q \in d_X | d) = 0$ (in particular, $P(q | d)$ need not be the OTS $P(q | d)$ discussed above). For such a scenario, for a vertical likelihood (i.e., a $P(d | f)$ that is independent of the values of $f(x \notin d_X, \cdot)$), we get the following result (see appendix C):

Theorem 2 For OTS error, a vertical $P(d | f)$, and a homogenous loss L , the uniform average over all targets f of $P(c | f, m) = \Lambda(c) / r$.

Again, this holds for any learning algorithm, and any sampling distribution. Note that this result means in particular that the "weight" of f 's on which one's algorithm performs worse than the random algorithm equals the weight for which it performs better. In other words, one can just as readily have a target for which one's algorithm has *worse than random guessing* as one in which it performs better than random. The pitfall we wish to avoid in supervised learning is not simply that our algorithm performs as poorly as random

We start with the following simple lemma, that recurs frequently in the subsequent analysis. Its proof is in appendix C.

Lemma 1 $P(c \mid f, d) = \sum_{y_H, y_F, q} \delta(c, L(y_H, y_F)) P(y_H \mid q, d) P(y_F \mid q, f) P(q \mid d)$.

Consider now the “(uniform) random learning algorithm”: for any test set element not in the training set, guess the output randomly (independently of the training set d), according to a uniform distribution. (With certain extra stipulations concerning behavior for test set questions $q \in d_X$, this is a version of the Gibbs learning algorithm.) An immediate implication of lemma (1), proven in appendix C, is that for this algorithm, for a symmetric homogenous loss function, $P(c \mid d) = \Lambda(c) / r$ for all training sets d . Similarly, for all priors over targets f , indicated by α , both $P(c \mid m, \alpha)$ and $P(c \mid d, \alpha)$ equal $\Lambda(c) / r$, for this random learning algorithm.

This simple kind of reasoning suffices to get “NFL” results for the random algorithm, even without invoking a vertical likelihood. However more is needed for scenarios concerning other algorithms, scenarios in which there is “randomness”, but it concerns targets rather than hypotheses. This is because we’re interested in probability distributions conditioned on target-based quantities (f, α , etc.), so results for when there is randomness in hypothesis-based quantities do not immediately carry over to results for randomness in target-based quantities.

To analyze these alternative scenarios, we start with the following simple implication of lemma (1) (see appendix C):

The uniform average over all targets f of $P(c \mid f, d)$ equals

$$(1 / r) \sum_{y_H, y_F, q} \delta(c, L(y_H, y_F)) P(y_H \mid q, d) P(q \mid d).$$

Recalling the definition of homogenous loss L , we have now proven the following:

Theorem 1 For homogenous loss L , the uniform average over all f of $P(c \mid f, d)$ equals $\Lambda(c) / r$.

Note that this f -average is independent of the learning algorithm. So theorem (1) constitutes an NFL theorem for distributions conditioned on targets f and training sets d ; it says that uniformly averaged over all f , such distributions are independent of the learning algorithm. Note that this result holds independent of the sampling distribution, the training set, or the likelihood.

$P(f | d)$ and $P(h | d)$, respectively.) Metaphorically speaking, how “aligned” you (the learning algorithm) are with the universe (the posterior) determines how well you will generalize.

The question arises though of how much can be said concerning a particular learning algorithm’s generalization behavior without specifying the posterior (which usually means without specifying the prior). More precisely, the goal is to address the issue of how F_1 , the set of targets f for which algorithm A outperforms algorithm B, compares to F_2 , the set of targets f for which the reverse is true. To analyze this issue, the simple trick is used of comparing the average over f of f -conditioned probability distributions for algorithm A to the same average for algorithm B. The relationship between those averages is then used to compare F_1 to F_2 .

Evaluating such f -averages results in a set of NFL theorems. In this section, first I derive the NFL theorems for the case where the target f need not be single-valued. In this case, the theorems say that uniformly averaged over all f , all learning algorithms are identical. The implications of this for how F_1 compares to F_2 are discussed after the derivation of the theorems.

When the target f is not single-valued, it is a (countable) set of real numbers (one for each possible x - y pair). Accordingly, $P(f)$ is a probability density function in a multi-dimensional space. That makes integrating over all $P(f)$ ’s a subtle mathematical exercise. However in the function+noise scenario, for a fixed noise process, “ f ” is indexed by a single-valued function ϕ . Since there are a countable number of ϕ ’s, $P(\phi)$ is a countable set of real numbers, and it is straight-forward to integrate over all $P(\phi)$. Doing so gives some more NFL theorems, where one uniformly averages over all priors rather than just over all targets f . These additional theorems are presented after those involving averages over all targets f .

After deriving these theorems, I present some examples of them, designed to highlight their counter-intuitive aspects. I also present a general discussion of the significance of the theorems, and in particular of the uniform averaging that goes into deriving them.

Here and throughout this paper, when discussing non-single-valued f ’s, “ $A(f)$ uniformly averaged over all targets f ” means $\int df A(f) / \int df 1$. Note that these integrals are implicitly restricted to those f that constitute \mathbf{X} -conditioned distributions over \mathbf{Y} , i.e., to the appropriate product space of unit-simplices. (The details will not matter, because integrals will almost never need to be evaluated. But formally, integrals over targets f are over a full r^n -dimensional Euclidean space, with a product of Dirac delta functions and Heaviside functions inside the integrand enforcing the restriction to the Cartesian product of simplices.)

Similar meanings for “uniformly averaged” are assumed if we’re talking about averaging over other quantities, like $P(\phi)$.

i) Averaging over all target distributions f

The sets \mathbf{X} and \mathbf{Y} , of sizes n and r :	The input and output space, respectively.
The set d , of m \mathbf{X} - \mathbf{Y} pairs:	The training set.
The \mathbf{X} -conditioned distribution over \mathbf{Y} , f :	The target, used to generate test sets.
The \mathbf{X} -conditioned distribution over \mathbf{Y} , h :	The hypothesis, used to guess for test sets.
The real number c :	The cost.
The \mathbf{X} -value q :	The test set point.
The \mathbf{Y} -value y_F :	The sample of the target f at point q .
The \mathbf{Y} -value y_H :	The sample of the hypothesis h at point q .
$P(h d)$:	The learning algorithm.
$P(f d)$:	The posterior.
$P(d f)$:	The likelihood.
$P(f)$:	The prior.
If $c = L(y_F, y_H)$, $L(., .)$ is the “ <u>loss function</u> ”.	
L is “ <u>homogenous</u> ” if $\sum_{y_F} \delta(c, L(y_H, y_F))$ is independent of y_F .	
If we restrict attention to f ’s given by a fixed noise process superimposed on an underlying single-valued function from \mathbf{X} to \mathbf{Y} , ϕ , and if $\sum_{\phi} P(y_F q, \phi)$ is independent of y_F , we have “ <u>homogenous</u> ” noise.	
<hr/>	
Table 1: Summary of the terms in the EBF.	

3. THE NO-FREE-LUNCH THEOREMS

In [Wolpert 1992] it is shown that $P(c | d) = \int df dh P(h | d) P(f | d) M_{c,d}(f, h)$, where so long as the loss function is symmetric in its arguments, $M_{c,d}(., .)$ is symmetric in its arguments. (See point (11) of the previous section.) In other words, for the most common kinds of loss functions (zero-one, quadratic, etc.), the probability of a particular cost is determined by an inner product between your learning algorithm and the posterior probability. (f and h being the component labels of the d -indexed infinite-dimensional vectors

loss function $L(y_H, y_F)$. As an example, in regression, often we have “quadratic loss”: $L(y_H, y_F) = (y_H - y_F)^2$.

$L(., .)$ is “homogenous” if the sum over y_F of $\delta(c, L(y_H, y_F))$ is some function $\Lambda(c)$, independent of y_H (δ here being the Kronecker delta function). As an example, the “zero-one” loss traditional in computational learning theory ($L(a, b) = 1$ if $a \neq b$, 0 otherwise) is homogenous.

13) In the case of “IID error” (the conventional error measure), $P(q | d) = \pi(q)$ (so test set inputs are chosen according to the same distribution that determines training set inputs). In the case of OTS error, $P(q | d) = [\delta(q \notin d_X) \pi(q)] / [\sum_q \delta(q \notin d_X) \pi(q)]$, where $\delta(z) \equiv 1$ if z is true, 0 otherwise.

Subscripts $_{OTS}$ or $_{IID}$ on c correspond to using those respective kinds of error.

14) The “generalization error function” used in much of supervised learning is given by $c' \equiv E(C | f, h, d)$. (Subscripts $_{OTS}$ or $_{IID}$ on c' correspond to using those respective ways to generate q .) It is the average over all q of the cost c , for a given target f , hypothesis h , and training set d .

In general, probability distributions over c' do not by themselves determine those over c nor vice-versa, i.e., there is not an injection between such distributions. However the results in this paper in general hold for both c and c' , although they will only be presented for c . In addition, especially when relating results in this paper to theorems in the literature, sometimes results for c' will implicitly be meant even when the text still refers to c . (The context will make this clear.)

15) When the size of \mathbf{X} , n , is much greater than the size of the training set, m , probability distributions over c'_{IID} and distributions over c'_{OTS} become identical. (Although as mentioned in the previous section, distributions conditioned on c'_{IID} can be drastically different from those conditioned on c'_{OTS} .) This is established formally in appendix B.

tical.

As another example, if there is noise in generating training set \mathbf{X} values but none for test set \mathbf{X} values, then we usually do not have a vertical $P(d | f)$. (This is because, formally speaking, f directly governs the generation of test sets, not training sets; see appendix A.)

7) The “posterior” usually means $P(f | d)$, and the “prior” usually means $P(f)$.

8) It will be convenient at times to restrict attention to f ’s that are constructed by adding noise to a single-valued function from \mathbf{X} to \mathbf{Y} , ϕ . For a fixed noise process, such f ’s are indexed by the underlying ϕ .

The noise process is “homogenous” if the sum over all ϕ of $P(y_F | q, \phi)$ is independent of y_F . An example of a homogenous noise process is classification noise that with probability p replaces $\phi(q)$ with some other value in \mathbf{Y} , where that “other value in \mathbf{Y} ” is chosen uniformly and randomly.

iii) The learning algorithm

9) Hypotheses h are always assumed to be of the form of \mathbf{X} -conditioned distributions over \mathbf{Y} , indicated by the real-valued function $h(x \in \mathbf{X}, y \in \mathbf{Y})$ (i.e., $P(y_H | h, q) = h(q, y_H)$). Equivalently, where S_r is defined as the r -dimensional unit simplex, hypotheses can be viewed as mappings $h: \mathbf{X} \rightarrow S_r$.

Any restrictions on h are imposed by $P(f, h, d, c)$. Here and throughout, a “single-valued” distribution is one that, for a given x , is a delta function about some y . Such a distribution is a single-valued function from \mathbf{X} to \mathbf{Y} . As an example, if one is using a neural net to as one’s regression through the training set, usually the (neural net) h is single-valued. On the other hand, when one is performing probabilistic classification (as in softmax), h isn’t single-valued.

10) Any (!) learning algorithm (aka “generalizer”) is given by $P(h | d)$, although writing down a learning algorithm’s $P(h | d)$ explicitly is often quite difficult. A learning algorithm is “deterministic” if the same d always gives the same h . Backprop with a random initial weight is not deterministic. Nearest neighbor is.

Note that since d is ordered, “on-line” learning algorithms are subsumed as a special case.

11) The learning algorithm only sees the training set d , and in particular does not directly see the target. So $P(h | f, d) = P(h | d)$, which means that $P(h, f | d) = P(h | d) \times P(f | d)$, and therefore $P(f | h, d) = P(h, f | d) / P(h | d) = P(f | d)$.

iv) The cost and “generalization error”

12) For the purposes of this paper, the cost c is associated with a particular y_H and y_F , and is given by a

licated by H), the target (i.e., “true”) \mathbf{X} - \mathbf{Y} relationship (F), the training set (D), and the real world cost (C).

These variables are related to one another through other random variables representing the (test set) input space value (Q), and the associated target and hypothesis \mathbf{Y} -values, Y_F and Y_H respectively (with instantiations y_F and y_H respectively).

This completes the list of random variables.

As an example of the relationship between these random variables and supervised learning, f , a particular instantiation of a target, could refer to a “teacher” neural net together with superimposed noise. This noise-corrupted neural net generates the training set d . The hypothesis h on the other hand could be the neural net made by one’s “student” algorithm after training on d . Then q would be an input element of the test set, y_F and y_H associated samples of the outputs of the two neural nets for that element (the sampling of y_F including the effects of the superimposed noise), and c the resultant “cost” (e.g., c could be $(y_F - y_H)^2$).

ii) Training sets and targets

4) m is the number of elements in the (ordered) training set d . $\{d_{\mathbf{X}}(i), d_{\mathbf{Y}}(i)\}$ is the set of m input and output values in d . m' is the number of distinct values in $d_{\mathbf{X}}$.

5) Targets f are always assumed to be of the form of \mathbf{X} -conditioned distributions over \mathbf{Y} , indicated by the real-valued function $f(x \in \mathbf{X}, y \in \mathbf{Y})$ (i.e., $P(y_F | f, q) = f(q, y_F)$). Equivalently, where S_r is defined as the r -dimensional unit simplex, targets can be viewed as mappings $f: \mathbf{X} \rightarrow S_r$.

Any restrictions on f are imposed by $P(f, h, d, c)$, and in particular by its marginalization, $P(f)$. Note that any output noise process is automatically reflected in $P(y_F | f, q)$. Note also that the equality $P(y_F | f, q) = f(q, y_F)$ only directly refers to the generation of test set elements; in general, training set elements can be generated from targets in a different manner.

6) The “likelihood” is $P(d | f)$. It says how d was generated from f . It is “vertical” if $P(d | f)$ is independent of the values $f(x, y_F)$ for those $x \notin d_{\mathbf{X}}$. As an example, the conventional IID likelihood is

$$(2.1) \quad P(d | f) = \prod_{i=1}^m \pi(d_{\mathbf{X}}(i)) f(d_{\mathbf{X}}(i), d_{\mathbf{Y}}(i))$$

(where $\pi(x)$ is the “sampling distribution”). In other words, under this likelihood d is created by repeatedly and independently choosing an input value $d_{\mathbf{X}}(i)$ by sampling $\pi(x)$, and then choosing an associated output value by sampling $f(d_{\mathbf{X}}(i), \cdot)$, the same distribution used to generate test set outputs. This likelihood is ver-

None of this means that the conventional IID error measure is “wrong”. No claim is being made that one “should not” test with the same process that generated the training set. Rather the claim is simply that OTS testing is an issue of major importance. In that it gives no credit for memorization, it is also the natural way to investigate whether one can make assumption-free statements concerning an algorithm’s generalization (!) ability.

2. THE EXTENDED BAYESIAN FORMALISM

These papers use the Extended Bayesian Formalism [Wolpert 1994a, Wolpert et al. 1995, Wolpert 1992]. In the current context, the EBF is just conventional probability theory, applied to the case where one has a different random variable for the hypothesis output by the learning algorithm and for the target relationship. It is this crucial distinction that separates the EBF from conventional Bayesian analysis, and that allows the EBF (unlike conventional Bayesian analysis) to subsume all other major mathematical treatments of supervised learning like computational learning theory, sampling theory statistics, etc. (See [Wolpert 1994a].)

This section presents a synopsis of the EBF. Points (2), (8), (14) and (15) below can be skipped in a first reading. A quick reference of this section’s synopsis can be found in Table 1.

Readers unsure of any aspects of this synopsis, and in particular unsure of any of the formal basis of the EBF or justifications for any of its assumptions, are directed to the detailed exposition of the EBF in appendix A.

i) Overview

1) The input and output spaces are \mathbf{X} and \mathbf{Y} , respectively. They contain n and r elements respectively. A generic element of \mathbf{X} is indicated by ‘ x ’, and a generic element of \mathbf{Y} is indicated by ‘ y ’.

2) Random variables are indicated using capital letters. Associated instantiations of a random variable are indicated using lower case letters. Note though that some quantities (e.g., the space \mathbf{X}) are neither random variables nor instantiations of random variables, and therefore their written case carries no significance.

Only rarely will it be necessary to refer to a random variable rather than an instantiation of it. In accord with standard statistics notation, “ $E(A | b)$ ” will be used to mean the expectation value of A given $B = b$, i.e., to mean $\int da a P(a | b)$. (Sums replace integrals if appropriate.)

3) The primary random variables are the hypothesis \mathbf{X} - \mathbf{Y} relationship output by the learning algorithm (in-

Bayes-optimal learning algorithm. See also the discussion in [Wolpert 1994a] concerning the statistical physics supervised learning formalism.

viii) Second, although it's usually true that a probability distribution over IID error will well-approximate the corresponding distribution over OTS error, distributions *conditioned* on IID error can differ drastically from distributions conditioned on OTS error. This can be very important in understanding the results of computational learning theory.

As an example of such a difference, let s be the empirical misclassification rate between a hypothesis and the target over the training set (i.e., the average number of disagreements over the training set), m the size of the training set, c'_{IID} the misclassification rate over all of the input space (the IID zero-one loss generalization error), and c'_{OTS} the misclassification rate over that part of the input space lying outside of the training set. (These terms are formally defined in the next section and at the beginning of section 4.) Assume a uniform sampling distribution over the input space, a uniform prior over target input-output relationships, and a noise-free IID likelihood governing the training set generation. Then $P(s \mid c'_{\text{IID}}, m)$, the probability of getting empirical misclassification rate s given global misclassification rate c'_{IID} , averaged over all training sets of size m , is just the binomial distribution $(c'_{\text{IID}})^{sm} (1 - c'_{\text{IID}})^{(m-sm)} C_{sm}^m$, where $C_b^a \equiv a! / [b! (a - b)!]$ (s can be viewed as the percentage of heads in m flips of a coin with bias c'_{IID} towards heads).

On the other hand, $P(s \mid c'_{\text{OTS}}, m)$, the probability of getting empirical misclassification rate s given off-training sets misclassification rate c'_{OTS} , averaged over all training sets of size m , is independent of c'_{OTS} . (This is proven in section 4 below.) So the dependence of the empirical misclassification rate on the global misclassification rate depends crucially on whether it is OTS or IID “global misclassification rate”.

ix) Third, often it is more straight-forward to calculate a certain quantity for OTS rather than IID error. In such cases, even if one's ultimate interest is IID error, it makes sense to instead calculate OTS error (assuming one is in a regime where OTS error well-approximates IID error).

As an example, OTS error results presented in section 4 below mean that when the training set is much smaller than the full input space, $P(c'_{\text{IID}} \mid s, m)$ is (arbitrarily close to) independent of s , if the prior over target input-output relationships is uniform. This holds despite VC results saying that independent of the prior, it is highly unlikely for c'_{IID} and s to differ significantly. (This may seem paradoxical at first. See the discussion in section 4 below for the “resolution”.)

The formal identity (in the appropriate limit) between a probability distribution over an OTS error and one over an IID error is established at the end of appendix B.

is often exclusively on generalizing to as yet unseen examples.

iv) In the real world, very often the process generating the training set is not the same as that governing testing. In such scenarios, the usual justification for testing with the same process that generated the training set (and with it the possibility that test sets overlap with training sets) doesn't apply.

One example of such a difference between testing and training is "active" or "query-based" or "membership-based" learning. In that kind of learning the learner chooses, perhaps dynamically, where in the input space the training set elements are to be. However, conventionally, there is no such control over the test set. So testing and training are governed by different processes.

As another example, say we wish to learn tertiary protein structure from primary structure and then use that to aid drug design. We *already know* what tertiary structure corresponds to the primary structures in the training set. So we will never have those structures in the "test set" (i.e., in the set of nucleotide sequences whose tertiary structure we wish to infer to aid the drug design process). We will *only* be interested in OTS error.

v) Distinguishing the regime where test examples coincide with the training set from the one where there is no overlap amounts to splitting supervised learning along its natural "cleavage plane". Since behavior can be radically different in the two regimes, it's hard to see why one wouldn't want to distinguish them.

vi) When the training set is much smaller than the full input space, the probability that a randomly chosen test set input value coincides with the training set is vanishingly small. So in such situations one expects the value of the OTS error to be well-approximated by the value of the conventional IID (independent identically distributed) error, an error which allows overlap between test sets and training sets.

One might suppose that in such a small training set regime there is no aspect of OTS error not addressable by instead calculating IID error. This is wrong though, as the following several points illustrate.

vii) First, even if OTS error is well approximated by IID error, it does not follow that quantities like the "derivatives" of the errors are close to one another. In particular, it does not follow that the sign of the slope of the learning curve - often an object of major interest - is the same for both errors over some region of interest.

As an example, in [Wolpert et al., 1995], it is shown that the expected OTS misclassification rate can *increase* with training set size, even if one averages both over training sets and targets, and even if one uses the Bayes-optimal learning algorithm. In contrast, it is also shown there that under those same conditions, the expected IID misclassification rate is strictly non-increasing as a function of training set size for the

The major extensions beyond this previous work that is contained in these two papers are: i) many more issues are analyzed (e.g., essentially all of paper two was not touched upon in the earlier work); and 2) many fewer restrictions are made (e.g., losses other than zero-one are considered, arbitrary kinds of noise are allowed, both hypotheses and targets are arbitrary probability distributions rather than single-valued functions from inputs to outputs, etc.)

1. OFF-TRAINING-SET ERROR

Many introductory supervised learning texts take the view that “the overall objective ... is to learn from samples and to generalize to new, *as yet unseen* cases” (italics mine - see [Weiss and Kulikowski 1991], for example). Similarly, in supervised learning it is common practice to try to avoid fitting the training set exactly, to try to avoid “overtraining”. One of the major rationales given for this is that if one overtrains, “the resulting (system) is unlikely to classify *additional points* (in the input space) correctly” (italics mine - see [Dietterich, 1990]). As another example, in [Blumer et al. 1987], we read that “the real value of a scientific explanation lies not in its ability to explain (what one has already seen), but in predicting events that have yet to (be seen)”. As a final example, in [Mitchell et al. 1994] we read that “(in Machine Learning we wish to know whether) any hypothesis found to approximate the target function well over a sufficiently large set of training examples will also approximate the target function well over *other unobserved* examples.”

This language makes clear that OTS behavior is a central concern of supervised learning, even though little theoretical work has been devoted to it to date. Some of the reasons for such concern are as follows.

i) In the low-noise (for outputs) regime, optimal behavior *on the training set* is trivially determined by look-up table memorization. Of course, this has nothing to do with behavior off of the training set; so in this regime, it is only such OTS behavior that is of interest.

ii) In particular, in that low-noise regime, if one uses a memorizing learning algorithm, then for test sets overlapping with training sets the upper limit of possible test set error values shrinks as the training set grows. If one doesn't correct for this when comparing behavior for different sizes of the training set (as when investigating learning curves), one is comparing apples and oranges. In that low-noise regime, correcting for this effect by renormalizing the range of possible error values is equivalent to requiring that test sets and training sets be distinct, i.e., is equivalent to using OTS error. (See [Wolpert 1994a].)

iii) In artificial intelligence - one of the primary fields concerned with supervised learning - the emphasis

Some of these alternative comparisons reveal no distinctions between algorithms, just like the comparisons in paper one. However some of the other alternative comparisons result in *a priori* distinctions between algorithms. In particular, it is pointed out in paper two that the equivalence of average between cross-validation and anti-cross-validation does not mean they have equivalent “head-to-head minimax” properties, and that algorithms can differ in those properties. Indeed, it may be that cross-validation has better head-to-head minimax properties than anti-cross-validation, and therefore can be *a priori* justified in that sense.

Of course, the analysis of paper one does not rule out the possibility that there are targets for which a particular learning algorithm works well compared to some other one. To address the non-trivial aspects of this issue, paper two discusses the case where one averages over hypotheses rather than targets. The results of such analyses hold for all possible priors, since they hold for all (fixed) targets. This allows them to be used to prove, as a particular example, that cross-validation cannot be justified as a Bayesian procedure. I.e., there is no prior over targets for which, *without regard for the learning algorithms in question*, one can conclude that one should choose between those algorithms based on minimal rather than (for example) maximal cross-validation error. In addition, it is noted that for a very natural restriction of the class of learning algorithms, one can distinguish between using minimal rather than maximal cross-validation error - and the result is that one should use maximal error (!).

All of the analysis up to this point assumes the loss function is in the same class as the zero-one loss function (which is assumed in almost all of computational learning theory). Paper two goes on to discuss other loss functions. In particular, the quadratic loss function modifies the preceding results considerably; for that loss function, there **are** algorithms that are *a priori* superior to other algorithms. However it is shown in paper two that no algorithm is superior to its “randomized” version, in which the set of potential fits to training sets is held fixed, but which fit is associated with which training set changes. In this sense one cannot *a priori* justify any particular learning algorithm, even for a quadratic loss function.

Finally, paper two ends with a brief overview of some open issues and discusses future work.

It cannot be emphasized enough that no claim whatsoever is being made in this first paper that all algorithms are equivalent *in practice*, in the real world. In particular, no claim is being made that one should not use cross-validation in the real world. (I have done so myself many times in the past and intend to do so again in the future.) The sole concern of this paper is what can (not) be formally inferred about the utility of various learning algorithms if one makes no assumptions concerning targets.

The work in these papers builds upon the analysis in [Wolpert 1992, 1993]. Some aspects of that analysis are nicely synopsized in [Schaffer 1993, 1994]. [Schaffer 1994] also contains an interesting discussion of the implications of the NFL theorems for real world learning, as does [Murphy and Pazzani 1994]. See also [Wolpert and Macready 1995] for related work in the field of combinatorial optimization.

of generalization error such overlap is allowed.) Section 1 of this first paper explains why such a measure of error is of interest, and in particular emphasizes that it is equivalent to (more conventional) IID error in many scenarios of interest. Those who already accept that OTS error is of interest can skip this section.

Section 2 presents the mathematical formalism used in this paper.

Section 3 presents the “no free lunch” (NFL) theorems (phrase due to D. Haussler). Some of those theorems show, loosely speaking, that for any two algorithms A and B, there are “as many” targets for which algorithm A has lower expected OTS error than algorithm B as vice-versa (whether one averages over training sets or not). In particular, such equivalence holds even if one of the algorithms is random guessing; there are “as many” targets for which any particular learning algorithm gets confused by the data and performs *worse than random* as for which it performs better. As another example of the NFL theorems, it is shown explicitly that A is equivalent to B when B is an algorithm that chooses between two hypotheses based on which *disagrees* more with the training set, and A is an algorithm that chooses based on which *agrees* more with the training set. Other NFL theorems are also derived, showing, for example, that there are as many priors over targets in which A beats B (i.e., has lower expected error than B) as vice-versa. In this, the quotes presented at the beginning of this section are misleading at best.

Next a set of simple examples is presented illustrating the theorems in scenarios in which their applicability is somewhat counter-intuitive. In particular, a brief discussion is presented of the fact that there are as many targets for which it is preferable to choose between two learning algorithms based on which has *larger* cross-validation error (“anti-cross-validation”) as based on which has smaller cross-validation error.

This section also contains the subsection “Extensions for non-uniform averaging” which extends the NFL results beyond uniform averages; as that subsection shows, one can, for example, only consider priors over targets that are highly structured, and it is still often true that all algorithms are equal. Also in this section is the subsection “On uniform averaging”, which provides the intellectual context for the analyses that result in the NFL theorems.

Section 4 discusses the NFL theorem’s implications for and relationship with computational learning theory. It starts with a discussion of empirical error and OTS error. This discussion makes clear that one must be very careful in trying to interpret uniform convergence (VC) results. In particular, it makes clear that one can **not** say: if empirical misclassification rate is low; the VC dimension of your generalizer is small; and the training set is large, then with high probability your OTS error is small. After this, the implications of the NFL results for active learning, and for “membership queries” algorithms and “punting” algorithms (those that may refuse to make a guess) are discussed.

Small and simple proofs of claims made in the text of this first paper are collected in appendix C.

Paper one concentrates on relative sizes of sets of targets and the associated senses in which all algorithms are *a priori* equivalent. In contrast, paper two concentrates on other ways to compare algorithms.

INTRODUCTION

Much of modern supervised learning theory gives the impression that one can deduce something about the efficacy of a particular learning algorithm (generalizer) without the need for any assumptions about the target input-output relationship one is trying to learn with that algorithm. At most, it would appear, to make such a deduction one has to know something about the training set as well as about the learning algorithm.

Consider for example the following quotes from some well-known papers: “Theoretical studies link the generalization error of a learning algorithm to the error on the training examples and the capacity of the learning algorithm (independent of concerns about the target)”;

“We have given bounds (independent of the target) on the training set size vs. neural net size needed such that valid generalization can be expected”;

“If our network can be trained to classify correctly ... $1 - (1 - \gamma) \epsilon$ of the k training examples, then the probability its [generalization] error is less than ϵ is at least [a function, independent of the target, of ϵ , γ , k , and the learning algorithm]”;

“There are algorithms that with high probability produce good approximators regardless of the target function ... We do not need to make any assumption about prior probabilities (of targets)”;

“To do Bayesian analysis, it is not necessary to work out the prior (over targets)”;

“This shows that (the probability distribution of generalization accuracy) gets concentrated at higher and higher accuracy values as more examples are learned (independent of the target).”

Similar statements can be found in the “proofs” that various supervised learning communities have offered for Occam’s razor [Blumer et al 1987, Berger and Jeffreys 1992; see also Wolpert 1994a, 1995]. There even exists a field (“agnostic learning” [Kearns et al. 1992]) whose expressed purpose is to create learning algorithms that are assuredly effective even in the absence of assumptions about the target.

Frequently the authors of these kinds of quotes understand that there are subtleties and caveats behind them. But the quotes taken at face value raise an intriguing question: can one actually get something for nothing in supervised learning? Can one get useful, caveat-free theoretical results that link the training set and the learning algorithm to generalization error, without making assumptions concerning the target? More generally, are there useful practical techniques that require no such assumptions? As a potential example of such a technique, note that people usually use cross-validation without making any assumptions about the underlying target, as though the technique were universally applicable.

This is the first of two papers that present an initial investigation of this issue. These papers can be viewed as an analysis of the mathematical “skeleton” of supervised learning, before the “flesh” of particular priors over targets are introduced. It should be emphasized that the work in these papers is very preliminary; even the “skeleton” of supervised learning is extremely rich and detailed. Much remains to be done.

The primary mathematical tool used in these papers is off-training set (OTS) generalization error, i.e., generalization error for test sets that contain no overlap with the training set. (In the conventional measure

THE LACK OF A *PRIORI* DISTINCTIONS BETWEEN LEARNING ALGORITHMS

by

David H. Wolpert

The Santa Fe Institute, 1399 Hyde Park Rd., Santa Fe, NM, 87501, USA (dhw@santafe.edu)

Abstract: This is the first of two papers that use off-training set (OTS) error to investigate the assumption-free relationship between learning algorithms. This first paper discusses the senses in which there are no *a priori* distinctions between learning algorithms. (The second paper discusses the senses in which there are such distinctions.) In this first paper it is shown, loosely speaking, that for any two algorithms A and B, there are “as many” targets (or priors over targets) for which A has lower expected OTS error than B as vice-versa, for loss functions like zero-one loss. In particular, this is true if A is cross-validation and B is “anti-cross-validation” (choose the learning algorithm with largest cross-validation error). This paper ends with a discussion of the implications of these results for computational learning theory. It is shown that one can **not** say: if empirical misclassification rate is low; the Vapnik-Chervonenkis dimension of your generalizer is small; and the training set is large, then with high probability your OTS error is small. Other implications for “membership queries” algorithms and “punting” algorithms are also discussed.

“Even after the observation of the frequent conjunction of objects, we have no reason to draw any inference concerning any object beyond those of which we have had experience.” - David Hume, in *A Treatise of Human Nature*, Book I, part 3, Section 12.